

*Wiley Series in Discrete Mathematics and Optimization*

AN  
**INTRODUCTION**  
TO  
**OPTIMIZATION**

**FOURTH EDITION**

Edwin K. P. Chong  
Stanislaw H. Żak

 **WILEY**

**WWW.**  
LINK AVAILABLE



# **AN INTRODUCTION TO OPTIMIZATION**

---

**WILEY SERIES IN  
DISCRETE MATHEMATICS AND OPTIMIZATION**

A complete list of titles in this series appears at the end of this volume.

---

---

# **AN INTRODUCTION TO OPTIMIZATION**

---

Fourth Edition

**Edwin K. P. Chong**

Colorado State University

**Stanislaw H. Zak**

Purdue University



**WILEY**

A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2013 by John Wiley & Sons, Inc. All rights reserved

Published by John Wiley & Sons, Inc., Hoboken, New Jersey

Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at [www.copyright.com](http://www.copyright.com). Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permission>.

**Limit of Liability/Disclaimer of Warranty:** While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at [www.wiley.com](http://www.wiley.com).

***Library of Congress Cataloging-in-Publication Data***

Chong, Edwin Kah Pin.

An introduction to optimization / Edwin K. P. Chong, Colorado State University, Stanislaw H. Zak, Purdue University. — Fourth edition.

pages cm

Summary: "The purpose of the book is to give the reader a working knowledge of optimization theory and methods" — Provided by publisher.

Includes bibliographical references and index.

ISBN 978-1-118-27901-4 (hardback)

1. Mathematical optimization. I. Zak, Stanislaw H. II. Title.

QA402.5.C476 2012

519.6—dc23

2012031772

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

*To my wife, Yat-Yee,  
and to my parents, Paul  
and Julianne Chong.  
Edwin K. P. Chong*

*To JMJ; my wife,  
Mary Ann; and my  
parents, Janina and  
Konstanty Żak.  
Stanislaw H. Żak*



# CONTENTS

---

Preface	xiii
---------	------

## PART I MATHEMATICAL REVIEW

<b>1 Methods of Proof and Some Notation</b>	<b>3</b>
1.1 Methods of Proof	3
1.2 Notation	5
Exercises	6
<b>2 Vector Spaces and Matrices</b>	<b>7</b>
2.1 Vector and Matrix	7
2.2 Rank of a Matrix	13
2.3 Linear Equations	17
2.4 Inner Products and Norms	19
Exercises	22
<b>3 Transformations</b>	<b>25</b>
3.1 Linear Transformations	25

3.2	Eigenvalues and Eigenvectors	26
3.3	Orthogonal Projections	29
3.4	Quadratic Forms	31
3.5	Matrix Norms	35
	Exercises	40
<b>4</b>	<b>Concepts from Geometry</b>	<b>45</b>
4.1	Line Segments	45
4.2	Hyperplanes and Linear Varieties	46
4.3	Convex Sets	48
4.4	Neighborhoods	50
4.5	Polytopes and Polyhedra	52
	Exercises	53
<b>5</b>	<b>Elements of Calculus</b>	<b>55</b>
5.1	Sequences and Limits	55
5.2	Differentiability	62
5.3	The Derivative Matrix	63
5.4	Differentiation Rules	67
5.5	Level Sets and Gradients	68
5.6	Taylor Series	72
	Exercises	77
<b>PART II UNCONSTRAINED OPTIMIZATION</b>		
<b>6</b>	<b>Basics of Set-Constrained and Unconstrained Optimization</b>	<b>81</b>
6.1	Introduction	81
6.2	Conditions for Local Minimizers	83
	Exercises	93
<b>7</b>	<b>One-Dimensional Search Methods</b>	<b>103</b>
7.1	Introduction	103
7.2	Golden Section Search	104
7.3	Fibonacci Method	108
7.4	Bisection Method	116
7.5	Newton's Method	116
7.6	Secant Method	120
7.7	Bracketing	123

7.8	Line Search in Multidimensional Optimization	124
	Exercises	126
<b>8</b>	<b>Gradient Methods</b>	<b>131</b>
8.1	Introduction	131
8.2	The Method of Steepest Descent	133
8.3	Analysis of Gradient Methods	141
	Exercises	153
<b>9</b>	<b>Newton's Method</b>	<b>161</b>
9.1	Introduction	161
9.2	Analysis of Newton's Method	164
9.3	Levenberg-Marquardt Modification	168
9.4	Newton's Method for Nonlinear Least Squares	168
	Exercises	171
<b>10</b>	<b>Conjugate Direction Methods</b>	<b>175</b>
10.1	Introduction	175
10.2	The Conjugate Direction Algorithm	177
10.3	The Conjugate Gradient Algorithm	182
10.4	The Conjugate Gradient Algorithm for Nonquadratic Problems	186
	Exercises	189
<b>11</b>	<b>Quasi-Newton Methods</b>	<b>193</b>
11.1	Introduction	193
11.2	Approximating the Inverse Hessian	194
11.3	The Rank One Correction Formula	197
11.4	The DFP Algorithm	202
11.5	The BFGS Algorithm	207
	Exercises	211
<b>12</b>	<b>Solving Linear Equations</b>	<b>217</b>
12.1	Least-Squares Analysis	217
12.2	The Recursive Least-Squares Algorithm	227
12.3	Solution to a Linear Equation with Minimum Norm	231
12.4	Kaczmarz's Algorithm	232
12.5	Solving Linear Equations in General	236

Exercises	244
<b>13 Unconstrained Optimization and Neural Networks</b>	<b>253</b>
13.1 Introduction	253
13.2 Single-Neuron Training	256
13.3 The Backpropagation Algorithm	258
Exercises	270
<b>14 Global Search Algorithms</b>	<b>273</b>
14.1 Introduction	273
14.2 The Nelder-Mead Simplex Algorithm	274
14.3 Simulated Annealing	278
14.4 Particle Swarm Optimization	282
14.5 Genetic Algorithms	285
Exercises	298
<b>PART III LINEAR PROGRAMMING</b>	
<b>15 Introduction to Linear Programming</b>	<b>305</b>
15.1 Brief History of Linear Programming	305
15.2 Simple Examples of Linear Programs	307
15.3 Two-Dimensional Linear Programs	314
15.4 Convex Polyhedra and Linear Programming	316
15.5 Standard Form Linear Programs	318
15.6 Basic Solutions	324
15.7 Properties of Basic Solutions	327
15.8 Geometric View of Linear Programs	330
Exercises	335
<b>16 Simplex Method</b>	<b>339</b>
16.1 Solving Linear Equations Using Row Operations	339
16.2 The Canonical Augmented Matrix	346
16.3 Updating the Augmented Matrix	349
16.4 The Simplex Algorithm	350
16.5 Matrix Form of the Simplex Method	357
16.6 Two-Phase Simplex Method	361
16.7 Revised Simplex Method	364
Exercises	369

<b>17 Duality</b>	<b>379</b>
17.1 Dual Linear Programs	379
17.2 Properties of Dual Problems	387
Exercises	394
<b>18 Nonsimplex Methods</b>	<b>403</b>
18.1 Introduction	403
18.2 Khachiyan's Method	405
18.3 Affine Scaling Method	408
18.4 Karmarkar's Method	413
Exercises	426
<b>19 Integer Linear Programming</b>	<b>429</b>
19.1 Introduction	429
19.2 Unimodular Matrices	430
19.3 The Gomory Cutting-Plane Method	437
Exercises	447
<b>PART IV NONLINEAR CONSTRAINED OPTIMIZATION</b>	
<b>20 Problems with Equality Constraints</b>	<b>453</b>
20.1 Introduction	453
20.2 Problem Formulation	455
20.3 Tangent and Normal Spaces	456
20.4 Lagrange Condition	463
20.5 Second-Order Conditions	472
20.6 Minimizing Quadratics Subject to Linear Constraints	476
Exercises	481
<b>21 Problems with Inequality Constraints</b>	<b>487</b>
21.1 Karush-Kuhn-Tucker Condition	487
21.2 Second-Order Conditions	496
Exercises	501
<b>22 Convex Optimization Problems</b>	<b>509</b>
22.1 Introduction	509
22.2 Convex Functions	512
22.3 Convex Optimization Problems	521

22.4	Semidefinite Programming	527
	Exercises	540
<b>23</b>	<b>Algorithms for Constrained Optimization</b>	<b>549</b>
23.1	Introduction	549
23.2	Projections	549
23.3	Projected Gradient Methods with Linear Constraints	553
23.4	Lagrangian Algorithms	557
23.5	Penalty Methods	564
	Exercises	571
<b>24</b>	<b>Multiobjective Optimization</b>	<b>577</b>
24.1	Introduction	577
24.2	Pareto Solutions	578
24.3	Computing the Pareto Front	581
24.4	From Multiobjective to Single-Objective Optimization	585
24.5	Uncertain Linear Programming Problems	588
	Exercises	596
<b>References</b>		<b>599</b>
<b>Index</b>		<b>609</b>

## PREFACE

---

Optimization is central to any problem involving decision making, whether in engineering or in economics. The task of decision making entails choosing among various alternatives. This choice is governed by our desire to make the “best” decision. The measure of goodness of the alternatives is described by an objective function or performance index. Optimization theory and methods deal with selecting the best alternative in the sense of the given objective function.

The area of optimization has received enormous attention in recent years, primarily because of the rapid progress in computer technology, including the development and availability of user-friendly software, high-speed and parallel processors, and artificial neural networks. A clear example of this phenomenon is the wide accessibility of optimization software tools such as the Optimization Toolbox of MATLAB<sup>1</sup> and the many other commercial software packages.

There are currently several excellent graduate textbooks on optimization theory and methods (e.g., [3], [39], [43], [51], [87], [88], [104], [129]), as well as undergraduate textbooks on the subject with an emphasis on engineering design (e.g., [1] and [109]). However, there is a need for an introductory

<sup>1</sup>MATLAB is a registered trademark of The MathWorks, Inc.

textbook on optimization theory and methods at a senior undergraduate or beginning graduate level. The present text was written with this goal in mind. The material is an outgrowth of our lecture notes for a one-semester course in optimization methods for seniors and beginning graduate students at Purdue University, West Lafayette, Indiana. In our presentation, we assume a working knowledge of basic linear algebra and multivariable calculus. For the reader's convenience, a part of this book (Part I) is devoted to a review of the required mathematical background material. Numerous figures throughout the text complement the written presentation of the material. We also include a variety of exercises at the end of each chapter. A solutions manual with complete solutions to the exercises is available from the publisher to instructors who adopt this text. Some of the exercises require using MATLAB. The student edition of MATLAB is sufficient for all of the MATLAB exercises included in the text. The MATLAB source listings for the MATLAB exercises are also included in the solutions manual.

The purpose of the book is to give the reader a working knowledge of optimization theory and methods. To accomplish this goal, we include many examples that illustrate the theory and algorithms discussed in the text. However, it is not our intention to provide a cookbook of the most recent numerical techniques for optimization; rather, our goal is to equip the reader with sufficient background for further study of advanced topics in optimization.

The field of optimization is still a very active research area. In recent years, various new approaches to optimization have been proposed. In this text, we have tried to reflect at least some of the flavor of recent activity in the area. For example, we include a discussion of randomized search methods—these include particle swarm optimization and genetic algorithms, topics of increasing importance in the study of complex adaptive systems. There has also been a recent surge of applications of optimization methods to a variety of new problems. An example of this is the use of descent algorithms for the training of feedforward neural networks. An entire chapter in the book is devoted to this topic. The area of neural networks is an active area of ongoing research, and many books have been devoted to this subject. The topic of neural network training fits perfectly into the framework of unconstrained optimization methods. Therefore, the chapter on feedforward neural networks not only provides an example of application of unconstrained optimization methods but also gives the reader an accessible introduction to what is currently a topic of wide interest.

The material in this book is organized into four parts. Part I contains a review of some basic definitions, notations, and relations from linear algebra, geometry, and calculus that we use frequently throughout the book. In Part II we consider unconstrained optimization problems. We first discuss some theoretical foundations of set-constrained and unconstrained optimization, including necessary and sufficient conditions for minimizers and maximizers. This is followed by a treatment of various iterative optimization algorithms, including line search methods, together with their properties. A discussion of global

search algorithms is included in this part. We also analyze the least-squares optimization problem and the associated recursive least-squares algorithm. Parts III and IV are devoted to constrained optimization. Part III deals with linear programming problems, which form an important class of constrained optimization problems. We give examples and analyze properties of linear programs, and then discuss the simplex method for solving linear programs. We also provide a brief treatment of dual linear programming problems. We then describe some nonsimplex algorithms for solving linear programs: Khachiyan's method, the affine scaling method, and Karmarkar's method. We wrap up Part III by discussing integer linear programming problems. In Part IV we treat nonlinear constrained optimization. Here, as in Part II, we first present some theoretical foundations of nonlinear constrained optimization problems, including convex optimization problems. We then discuss different algorithms for solving constrained optimization problems. We also treat multiobjective optimization.

Although we have made every effort to ensure an error-free text, we suspect that some errors remain undetected. For this purpose, we provide online updated errata that can be found at the Web site for the book, accessible via

<http://www.wiley.com/mathematics>

We are grateful to several people for their help during the course of writing this book. In particular, we thank Dennis Goodman of Lawrence Livermore Laboratories for his comments on early versions of Part II and for making available to us his lecture notes on nonlinear optimization. We thank Moshe Kam of Drexel University for pointing out some useful references on nonsimplex methods. We are grateful to Ed Silverman and Russell Quong for their valuable remarks on Part I of the first edition. We also thank the students of ECE 580 at Purdue University and ECE 520 and MATH 520 at Colorado State University for their many helpful comments and suggestions. In particular, we are grateful to Christopher Taylor for his diligent proofreading of early manuscripts of this book. This fourth edition incorporates many valuable suggestions of users of the first, second, and third editions, to whom we are grateful.

E. K. P. CHONG AND S. H. ŽAK  
*Fort Collins, Colorado, and West Lafayette, Indiana*



## **PART I**

---

# **MATHEMATICAL REVIEW**

---



# CHAPTER 1

---

## METHODS OF PROOF AND SOME NOTATION

---

### 1.1 Methods of Proof

Consider two statements, “A” and “B,” which could be either true or false. For example, let “A” be the statement “John is an engineering student,” and let “B” be the statement “John is taking a course on optimization.” We can combine these statements to form other statements, such as “A and B” or “A or B.” In our example, “A and B” means “John is an engineering student, and he is taking a course on optimization.” We can also form statements such as “not A,” “not B,” “not (A and B),” and so on. For example, “not A” means “John is not an engineering student.” The truth or falsity of the combined statements depend on the truth or falsity of the original statements, “A” and “B.” This relationship is expressed by means of truth tables; see Tables 1.1 and 1.2.

From the tables, it is easy to see that the statement “not (A and B)” is equivalent to “(not A) or (not B)” (see Exercise 1.3). This is called *DeMorgan’s law*.

In proving statements, it is convenient to express a combined statement by a *conditional*, such as “A implies B,” which we denote “ $A \Rightarrow B$ .” The conditional

**Table 1.1** Truth Table for “A and B” and “A or B”

A	B	A and B	A or B
F	F	F	F
F	T	F	T
T	F	F	T
T	T	T	T

**Table 1.2** Truth Table for “not A”

A	not A
F	T
T	F

**Table 1.3** Truth Table for Conditionals and Biconditionals

A	B	$A \Rightarrow B$	$A \Leftarrow B$	$A \Leftrightarrow B$
F	F	T	T	T
F	T	T	F	F
T	F	F	T	F
T	T	T	T	T

“ $A \Rightarrow B$ ” is simply the combined statement “(not A) or B” and is often also read “A only if B,” or “if A then B,” or “A is sufficient for B,” or “B is necessary for A.”

We can combine two conditional statements to form a *biconditional* statement of the form “ $A \Leftrightarrow B$ ,” which simply means “ $(A \Rightarrow B)$  and  $(B \Rightarrow A)$ .” The statement “ $A \Leftrightarrow B$ ” reads “A if and only if B,” or “A is equivalent to B,” or “A is necessary and sufficient for B.” Truth tables for conditional and biconditional statements are given in Table 1.3.

It is easy to verify, using the truth table, that the statement “ $A \Rightarrow B$ ” is equivalent to the statement “ $(\text{not } B) \Rightarrow (\text{not } A)$ .” The latter is called the *contrapositive* of the former. If we take the contrapositive to DeMorgan’s law, we obtain the assertion that “not (A or B)” is equivalent to “(not A) and (not B).”

Most statements we deal with have the form “ $A \Rightarrow B$ .” To prove such a statement, we may use one of the following three different techniques:

1. The direct method

2. Proof by contraposition
3. Proof by contradiction or *reductio ad absurdum*

In the case of the *direct method*, we start with “A,” then deduce a chain of various consequences to end with “B.”

A useful method for proving statements is *proof by contraposition*, based on the equivalence of the statements “ $A \Rightarrow B$ ” and “ $(\text{not } B) \Rightarrow (\text{not } A)$ .” We start with “not B,” then deduce various consequences to end with “not A” as a conclusion.

Another method of proof that we use is *proof by contradiction*, based on the equivalence of the statements “ $A \Rightarrow B$ ” and “not (A and (not B)).” Here we begin with “A and (not B)” and derive a contradiction.

Occasionally, we use the *principle of induction* to prove statements. This principle may be stated as follows. Assume that a given property of positive integers satisfies the following conditions:

- The number 1 possesses this property.
- If the number  $n$  possesses this property, then the number  $n + 1$  possesses it too.

The principle of induction states that under these assumptions any positive integer possesses the property.

The principle of induction is easily understood using the following intuitive argument. If the number 1 possesses the given property, then the second condition implies that the number 2 possesses the property. But, then again, the second condition implies that the number 3 possesses this property, and so on. The principle of induction is a formal statement of this intuitive reasoning.

For a detailed treatment of different methods of proof, see [130].

## 1.2 Notation

Throughout, we use the following notation. If  $X$  is a set, then we write  $x \in X$  to mean that  $x$  is an element of  $X$ . When an object  $x$  is not an element of a set  $X$ , we write  $x \notin X$ . We also use the “curly bracket notation” for sets, writing down the first few elements of a set followed by three dots. For example,  $\{x_1, x_2, x_3, \dots\}$  is the set containing the elements  $x_1, x_2, x_3$ , and so on. Alternatively, we can explicitly display the law of formation. For example,  $\{x : x \in \mathbb{R}, x > 5\}$  reads “the set of all  $x$  such that  $x$  is real and  $x$  is greater than 5.” The colon following  $x$  reads “such that.” An alternative notation for the same set is  $\{x \in \mathbb{R} : x > 5\}$ .

If  $X$  and  $Y$  are sets, then we write  $X \subset Y$  to mean that every element of  $X$  is also an element of  $Y$ . In this case, we say that  $X$  is a *subset* of  $Y$ . If  $X$  and  $Y$  are sets, then we denote by  $X \setminus Y$  (“ $X$  minus  $Y$ ”) the set of all points in  $X$  that are not in  $Y$ . Note that  $X \setminus Y$  is a subset of  $X$ . The

notation  $f : X \rightarrow Y$  means “ $f$  is a function from the set  $X$  into the set  $Y$ .” The symbol  $:=$  denotes arithmetic assignment. Thus, a statement of the form  $x := y$  means “ $x$  becomes  $y$ .” The symbol  $\triangleq$  means “equals by definition.”

Throughout the text, we mark the end of theorems, lemmas, propositions, and corollaries using the symbol  $\square$ . We mark the end of proofs, definitions, and examples by  $\blacksquare$ .

We use the IEEE style when citing reference items. For example, [77] represents reference number 77 in the list of references at the end of the book.

## EXERCISES

**1.1** Construct the truth table for the statement “ $(\text{not } B) \Rightarrow (\text{not } A)$ ,” and use it to show that this statement is equivalent to the statement “ $A \Rightarrow B$ .”

**1.2** Construct the truth table for the statement “ $\text{not } (A \text{ and } (\text{not } B))$ ,” and use it to show that this statement is equivalent to the statement “ $A \Rightarrow B$ .”

**1.3** Prove DeMorgan’s law by constructing the appropriate truth tables.

**1.4** Prove that for any statements  $A$  and  $B$ , we have “ $A \Leftrightarrow (A \text{ and } B)$  or  $(A \text{ and } (\text{not } B))$ .” This is useful because it allows us to prove a statement  $A$  by proving the two separate cases “ $(A \text{ and } B)$ ” and “ $(A \text{ and } (\text{not } B))$ .” For example, to prove that  $|x| \geq x$  for any  $x \in \mathbb{R}$ , we separately prove the cases “ $|x| \geq x$  and  $x \geq 0$ ” and “ $|x| \geq x$  and  $x < 0$ .” Proving the two cases turns out to be easier than proving the statement  $|x| \geq x$  directly (see Section 2.4 and Exercise 2.7).

**1.5** (This exercise is adopted from [22, pp. 80–81]) Suppose that you are shown four cards, laid out in a row. Each card has a letter on one side and a number on the other. On the visible side of the cards are printed the symbols

$S \quad 8 \quad 3 \quad A$

Determine which cards you should turn over to decide if the following rule is true or false: “If there is a vowel on one side of the card, then there is an even number on the other side.”

## CHAPTER 2

---

# VECTOR SPACES AND MATRICES

---

### 2.1 Vector and Matrix

We define a *column n*-vector to be an array of  $n$  numbers, denoted

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}.$$

The number  $a_i$  is called the  $i$ th component of the vector  $\mathbf{a}$ . Denote by  $\mathbb{R}$  the set of real numbers and by  $\mathbb{R}^n$  the set of column  $n$ -vectors with real components. We call  $\mathbb{R}^n$  an  $n$ -dimensional *real vector space*. We commonly denote elements of  $\mathbb{R}^n$  by lowercase bold letters (e.g.,  $\mathbf{x}$ ). The components of  $\mathbf{x} \in \mathbb{R}^n$  are denoted  $x_1, \dots, x_n$ .

We define a *row n*-vector as

$$[a_1, a_2, \dots, a_n].$$

The *transpose* of a given column vector  $\mathbf{a}$  is a row vector with corresponding elements, denoted  $\mathbf{a}^\top$ . For example, if

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix},$$

then

$$\mathbf{a}^\top = [a_1, a_2, \dots, a_n].$$

Equivalently, we may write  $\mathbf{a} = [a_1, a_2, \dots, a_n]^\top$ . Throughout the text we adopt the convention that the term *vector* (without the qualifier *row* or *column*) refers to a column vector.

Two vectors  $\mathbf{a} = [a_1, a_2, \dots, a_n]^\top$  and  $\mathbf{b} = [b_1, b_2, \dots, b_n]^\top$  are equal if  $a_i = b_i$ ,  $i = 1, 2, \dots, n$ .

The sum of the vectors  $\mathbf{a}$  and  $\mathbf{b}$ , denoted  $\mathbf{a} + \mathbf{b}$ , is the vector

$$\mathbf{a} + \mathbf{b} = [a_1 + b_1, a_2 + b_2, \dots, a_n + b_n]^\top.$$

The operation of addition of vectors has the following properties:

1. The operation is commutative:

$$\mathbf{a} + \mathbf{b} = \mathbf{b} + \mathbf{a}.$$

2. The operation is associative:

$$(\mathbf{a} + \mathbf{b}) + \mathbf{c} = \mathbf{a} + (\mathbf{b} + \mathbf{c}).$$

3. There is a zero vector

$$\mathbf{0} = [0, 0, \dots, 0]^\top$$

such that

$$\mathbf{a} + \mathbf{0} = \mathbf{0} + \mathbf{a} = \mathbf{a}.$$

The vector

$$[a_1 - b_1, a_2 - b_2, \dots, a_n - b_n]^\top$$

is called the difference between  $\mathbf{a}$  and  $\mathbf{b}$  and is denoted  $\mathbf{a} - \mathbf{b}$ .

The vector  $\mathbf{0} - \mathbf{b}$  is denoted  $-\mathbf{b}$ . Note that

$$\begin{aligned} \mathbf{b} + (\mathbf{a} - \mathbf{b}) &= \mathbf{a}, \\ -(-\mathbf{b}) &= \mathbf{b}, \\ -(\mathbf{a} - \mathbf{b}) &= \mathbf{b} - \mathbf{a}. \end{aligned}$$

The vector  $\mathbf{b} - \mathbf{a}$  is the unique solution of the vector equation

$$\mathbf{a} + \mathbf{x} = \mathbf{b}.$$

Indeed, suppose that  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  is a solution to  $\mathbf{a} + \mathbf{x} = \mathbf{b}$ . Then,

$$\begin{aligned} a_1 + x_1 &= b_1, \\ a_2 + x_2 &= b_2, \\ &\vdots \\ a_n + x_n &= b_n, \end{aligned}$$

and thus

$$\mathbf{x} = \mathbf{b} - \mathbf{a}.$$

We define an operation of multiplication of a vector  $\mathbf{a} \in \mathbb{R}^n$  by a real scalar  $\alpha \in \mathbb{R}$  as

$$\alpha\mathbf{a} = [\alpha a_1, \alpha a_2, \dots, \alpha a_n]^\top.$$

This operation has the following properties:

1. The operation is distributive: for any real scalars  $\alpha$  and  $\beta$ ,

$$\begin{aligned} \alpha(\mathbf{a} + \mathbf{b}) &= \alpha\mathbf{a} + \alpha\mathbf{b}, \\ (\alpha + \beta)\mathbf{a} &= \alpha\mathbf{a} + \beta\mathbf{a}. \end{aligned}$$

2. The operation is associative:

$$\alpha(\beta\mathbf{a}) = (\alpha\beta)\mathbf{a}.$$

3. The scalar 1 satisfies

$$1\mathbf{a} = \mathbf{a}.$$

4. Any scalar  $\alpha$  satisfies

$$\alpha\mathbf{0} = \mathbf{0}.$$

5. The scalar 0 satisfies

$$0\mathbf{a} = \mathbf{0}.$$

6. The scalar  $-1$  satisfies

$$(-1)\mathbf{a} = -\mathbf{a}.$$

Note that  $\alpha\mathbf{a} = \mathbf{0}$  if and only if  $\alpha = 0$  or  $\mathbf{a} = \mathbf{0}$ . To see this, observe that  $\alpha\mathbf{a} = \mathbf{0}$  is equivalent to  $\alpha a_1 = \alpha a_2 = \dots = \alpha a_n = 0$ . If  $\alpha = 0$  or  $\mathbf{a} = \mathbf{0}$ , then  $\alpha\mathbf{a} = \mathbf{0}$ . If  $\mathbf{a} \neq \mathbf{0}$ , then at least one of its components  $a_k \neq 0$ . For this component,  $\alpha a_k = 0$ , and hence we must have  $\alpha = 0$ . Similar arguments can be applied to the case when  $\alpha \neq 0$ .

A set of vectors  $\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$  is said to be *linearly independent* if the equality

$$\alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k = \mathbf{0}$$

implies that all coefficients  $\alpha_i$ ,  $i = 1, \dots, k$ , are equal to zero. A set of the vectors  $\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$  is *linearly dependent* if it is not linearly independent.

Note that the set composed of the single vector  $\mathbf{0}$  is linearly dependent, for if  $\alpha \neq 0$ , then  $\alpha \mathbf{0} = \mathbf{0}$ . In fact, any set of vectors containing the vector  $\mathbf{0}$  is linearly dependent.

A set composed of a single nonzero vector  $\mathbf{a} \neq \mathbf{0}$  is linearly independent since  $\alpha \mathbf{a} = \mathbf{0}$  implies that  $\alpha = 0$ .

A vector  $\mathbf{a}$  is said to be a *linear combination* of vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$  if there are scalars  $\alpha_1, \dots, \alpha_k$  such that

$$\mathbf{a} = \alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k.$$

**Proposition 2.1** *A set of vectors  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$  is linearly dependent if and only if one of the vectors from the set is a linear combination of the remaining vectors.*  $\square$

*Proof.*  $\Rightarrow$ : If  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$  is linearly dependent, then

$$\alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k = \mathbf{0},$$

where at least one of the scalars  $\alpha_i \neq 0$ , whence

$$\mathbf{a}_i = -\frac{\alpha_1}{\alpha_i} \mathbf{a}_1 - \frac{\alpha_2}{\alpha_i} \mathbf{a}_2 - \cdots - \frac{\alpha_k}{\alpha_i} \mathbf{a}_k.$$

$\Leftarrow$ : Suppose that

$$\mathbf{a}_1 = \alpha_2 \mathbf{a}_2 + \alpha_3 \mathbf{a}_3 + \cdots + \alpha_k \mathbf{a}_k,$$

then

$$(-1)\mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k = \mathbf{0}.$$

Because the first scalar is nonzero, the set of vectors  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$  is linearly dependent. The same argument holds if  $\mathbf{a}_i$ ,  $i = 2, \dots, k$ , is a linear combination of the remaining vectors.  $\blacksquare$

A subset  $\mathcal{V}$  of  $\mathbb{R}^n$  is called a *subspace* of  $\mathbb{R}^n$  if  $\mathcal{V}$  is closed under the operations of vector addition and scalar multiplication. That is, if  $\mathbf{a}$  and  $\mathbf{b}$  are vectors in  $\mathcal{V}$ , then the vectors  $\mathbf{a} + \mathbf{b}$  and  $\alpha \mathbf{a}$  are also in  $\mathcal{V}$  for every scalar  $\alpha$ .

Every subspace contains the zero vector  $\mathbf{0}$ , for if  $\mathbf{a}$  is an element of the subspace, so is  $(-1)\mathbf{a} = -\mathbf{a}$ . Hence,  $\mathbf{a} - \mathbf{a} = \mathbf{0}$  also belongs to the subspace.

Let  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$  be arbitrary vectors in  $\mathbb{R}^n$ . The set of all their linear combinations is called the *span* of  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$  and is denoted

$$\text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k] = \left\{ \sum_{i=1}^k \alpha_i \mathbf{a}_i : \alpha_1, \dots, \alpha_k \in \mathbb{R} \right\}.$$

Given a vector  $\mathbf{a}$ , the subspace  $\text{span}[\mathbf{a}]$  is composed of the vectors  $\alpha\mathbf{a}$ , where  $\alpha$  is an arbitrary real number ( $\alpha \in \mathbb{R}$ ). Also observe that if  $\mathbf{a}$  is a linear combination of  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$ , then

$$\text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k, \mathbf{a}] = \text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k].$$

The span of any set of vectors is a subspace.

Given a subspace  $\mathcal{V}$ , any set of linearly independent vectors  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\} \subset \mathcal{V}$  such that  $\mathcal{V} = \text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k]$  is referred to as a *basis* of the subspace  $\mathcal{V}$ . All bases of a subspace  $\mathcal{V}$  contain the same number of vectors. This number is called the *dimension* of  $\mathcal{V}$ , denoted  $\dim \mathcal{V}$ .

**Proposition 2.2** *If  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$  is a basis of  $\mathcal{V}$ , then any vector  $\mathbf{a}$  of  $\mathcal{V}$  can be represented uniquely as*

$$\mathbf{a} = \alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k,$$

where  $\alpha_i \in \mathbb{R}$ ,  $i = 1, 2, \dots, k$ . □

*Proof.* To prove the uniqueness of the representation of  $\mathbf{a}$  in terms of the basis vectors, assume that

$$\mathbf{a} = \alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k$$

and

$$\mathbf{a} = \beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \cdots + \beta_k \mathbf{a}_k.$$

We now show that  $\alpha_i = \beta_i$ ,  $i = 1, \dots, k$ . We have

$$\alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k = \beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \cdots + \beta_k \mathbf{a}_k$$

or

$$(\alpha_1 - \beta_1) \mathbf{a}_1 + (\alpha_2 - \beta_2) \mathbf{a}_2 + \cdots + (\alpha_k - \beta_k) \mathbf{a}_k = \mathbf{0}.$$

Because the set  $\{\mathbf{a}_i : i = 1, 2, \dots, k\}$  is linearly independent,  $\alpha_1 - \beta_1 = \alpha_2 - \beta_2 = \cdots = \alpha_k - \beta_k = 0$ , which implies that  $\alpha_i = \beta_i$ ,  $i = 1, \dots, k$ . ■

Suppose that we are given a basis  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$  of  $\mathcal{V}$  and a vector  $\mathbf{a} \in \mathcal{V}$  such that

$$\mathbf{a} = \alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \cdots + \alpha_k \mathbf{a}_k.$$

The coefficients  $\alpha_i$ ,  $i = 1, \dots, k$ , are called the *coordinates* of  $\mathbf{a}$  with respect to the basis  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}$ .

The *natural basis* for  $\mathbb{R}^n$  is the set of vectors

$$\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \quad \dots, \quad \mathbf{e}_n = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}.$$

The reason for calling these vectors the natural basis is that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \cdots + x_n \mathbf{e}_n.$$

We can similarly define *complex vector spaces*. For this, let  $\mathbb{C}$  denote the set of complex numbers and  $\mathbb{C}^n$  the set of column  $n$ -vectors with complex components. As the reader can easily verify, the set  $\mathbb{C}^n$  has properties similar to those of  $\mathbb{R}^n$ , where scalars can take complex values.

A *matrix* is a rectangular array of numbers, commonly denoted by uppercase bold letters (e.g.,  $\mathbf{A}$ ). A matrix with  $m$  rows and  $n$  columns is called an  $m \times n$  matrix, and we write

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

The real number  $a_{ij}$  located in the  $i$ th row and  $j$ th column is called the  $(i, j)$ th *entry*. We can think of  $\mathbf{A}$  in terms of its  $n$  columns, each of which is a column vector in  $\mathbb{R}^m$ . Alternatively, we can think of  $\mathbf{A}$  in terms of its  $m$  rows, each of which is a row  $n$ -vector.

Consider the  $m \times n$  matrix  $\mathbf{A}$  above. The *transpose* of matrix  $\mathbf{A}$ , denoted  $\mathbf{A}^\top$ , is the  $n \times m$  matrix

$$\mathbf{A}^\top = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix};$$

that is, the columns of  $\mathbf{A}$  are the rows of  $\mathbf{A}^\top$ , and vice versa.

Let the symbol  $\mathbb{R}^{m \times n}$  denote the set of  $m \times n$  matrices whose entries are real numbers. We treat column vectors in  $\mathbb{R}^n$  as elements of  $\mathbb{R}^{n \times 1}$ . Similarly, we treat row  $n$ -vectors as elements of  $\mathbb{R}^{1 \times n}$ . Accordingly, vector transposition is simply a special case of matrix transposition, and we will no longer distinguish between the two. Note that there is a slight inconsistency in the notation of row vectors when identified as  $1 \times n$  matrices: We separate the components of the row vector with commas, whereas in matrix notation we do not generally use commas. However, the use of commas in separating elements in a row helps to clarify their separation. We use such commas even in separating matrices arranged in a horizontal row.

## 2.2 Rank of a Matrix

Consider the  $m \times n$  matrix

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

Let us denote the  $k$ th column of  $\mathbf{A}$  by  $\mathbf{a}_k$ :

$$\mathbf{a}_k = \begin{bmatrix} a_{1k} \\ a_{2k} \\ \vdots \\ a_{mk} \end{bmatrix}.$$

The maximal number of linearly independent columns of  $\mathbf{A}$  is called the *rank* of the matrix  $\mathbf{A}$ , denoted  $\text{rank } \mathbf{A}$ . Note that  $\text{rank } \mathbf{A}$  is the dimension of  $\text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n]$ .

**Proposition 2.3** *The rank of a matrix  $\mathbf{A}$  is invariant under the following operations:*

1. Multiplication of the columns of  $\mathbf{A}$  by nonzero scalars.
2. Interchange of the columns.
3. Addition to a given column a linear combination of other columns.  $\square$

*Proof.*

1. Let  $\mathbf{b}_k = \alpha_k \mathbf{a}_k$ , where  $\alpha_k \neq 0$ ,  $k = 1, \dots, n$ , and let  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n]$ . Obviously,

$$\text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] = \text{span}[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n],$$

and thus

$$\text{rank } \mathbf{A} = \text{rank } \mathbf{B}.$$

2. The number of linearly independent vectors does not depend on their order.
3. Let

$$\mathbf{b}_1 = \mathbf{a}_1 + c_2 \mathbf{a}_2 + \cdots + c_n \mathbf{a}_n,$$

$$\mathbf{b}_2 = \mathbf{a}_2,$$

$\vdots$

$$\mathbf{b}_n = \mathbf{a}_n.$$

So, for any  $\alpha_1, \dots, \alpha_n$ ,

$$\alpha_1 \mathbf{b}_1 + \alpha_2 \mathbf{b}_2 + \cdots + \alpha_n \mathbf{b}_n = \alpha_1 \mathbf{a}_1 + (\alpha_2 + \alpha_1 c_2) \mathbf{a}_2 + \cdots + (\alpha_n + \alpha_1 c_n) \mathbf{a}_n,$$

and hence

$$\text{span}[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n] \subset \text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n].$$

On the other hand,

$$\begin{aligned}\mathbf{a}_1 &= \mathbf{b}_1 - c_2 \mathbf{b}_2 - \cdots - c_n \mathbf{b}_n, \\ \mathbf{a}_2 &= \mathbf{b}_2, \\ &\vdots \\ \mathbf{a}_n &= \mathbf{b}_n.\end{aligned}$$

Hence,

$$\text{span}[\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n] \subset \text{span}[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n].$$

Therefore,  $\text{rank } \mathbf{A} = \text{rank } \mathbf{B}$ . ■

A matrix  $\mathbf{A}$  is said to be *square* if the number of its rows is equal to the number of its columns (i.e., it is  $n \times n$ ). Associated with each square matrix  $\mathbf{A}$  is a scalar called the *determinant* of the matrix  $\mathbf{A}$ , denoted  $\det \mathbf{A}$  or  $|\mathbf{A}|$ . The determinant of a square matrix is a function of its columns and has the following properties:

1. The determinant of the matrix  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$  is a linear function of each column; that is,

$$\begin{aligned}\det[\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \alpha \mathbf{a}_k^{(1)} + \beta \mathbf{a}_k^{(2)}, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n] \\ = \alpha \det[\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \mathbf{a}_k^{(1)}, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n] \\ + \beta \det[\mathbf{a}_1, \dots, \mathbf{a}_{k-1}, \mathbf{a}_k^{(2)}, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n]\end{aligned}$$

for each  $\alpha, \beta \in \mathbb{R}$ ,  $\mathbf{a}_k^{(1)}, \mathbf{a}_k^{(2)} \in \mathbb{R}^n$ .

2. If for some  $k$  we have  $\mathbf{a}_k = \mathbf{a}_{k+1}$ , then

$$\det \mathbf{A} = \det[\mathbf{a}_1, \dots, \mathbf{a}_k, \mathbf{a}_{k+1}, \dots, \mathbf{a}_n] = \det[\mathbf{a}_1, \dots, \mathbf{a}_k, \mathbf{a}_k, \dots, \mathbf{a}_n] = 0.$$

3. Let

$$\mathbf{I}_n = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n] = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix},$$

where  $\{e_1, \dots, e_n\}$  is the natural basis for  $\mathbb{R}^n$ . Then

$$\det I_n = 1.$$

Note that if  $\alpha = \beta = 0$  in property 1, then

$$\det[a_1, \dots, a_{k-1}, 0, a_{k+1}, \dots, a_n] = 0.$$

Thus, if one of the columns is  $\mathbf{0}$ , then the determinant is equal to zero.

The determinant does not change its value if we add to a column another column multiplied by a scalar. This follows from properties 1 and 2 as shown below:

$$\begin{aligned} \det[a_1, \dots, a_{k-1}, a_k + \alpha a_j, a_{k+1}, \dots, a_j, \dots, a_n] \\ = \det[a_1, \dots, a_{k-1}, a_k, a_{k+1}, \dots, a_j, \dots, a_n] \\ + \alpha \det[a_1, \dots, a_{k-1}, a_j, a_{k+1}, \dots, a_j, \dots, a_n] \\ = \det[a_1, \dots, a_n]. \end{aligned}$$

However, the determinant changes its sign if we interchange columns. To show this property, note that

$$\begin{aligned} \det[a_1, \dots, a_{k-1}, a_k, a_{k+1}, \dots, a_n] \\ = \det[a_1, \dots, a_k + a_{k+1}, a_{k+1}, \dots, a_n] \\ = \det[a_1, \dots, a_k + a_{k+1}, a_{k+1} - (a_k + a_{k+1}), \dots, a_n] \\ = \det[a_1, \dots, a_k + a_{k+1}, -a_k, \dots, a_n] \\ = -\det[a_1, \dots, a_k + a_{k+1}, a_k, \dots, a_n] \\ = -(\det[a_1, \dots, a_k, a_k, \dots, a_n] + \det[a_1, \dots, a_{k+1}, a_k, \dots, a_n]) \\ = -\det[a_1, \dots, a_{k+1}, a_k, \dots, a_n]. \end{aligned}$$

A  $p$ th-order *minor* of an  $m \times n$  matrix  $\mathbf{A}$ , with  $p \leq \min\{m, n\}$ , is the determinant of a  $p \times p$  matrix obtained from  $\mathbf{A}$  by deleting  $m - p$  rows and  $n - p$  columns. (The notation  $\min\{m, n\}$  represents the smaller of  $m$  and  $n$ .)

We can use minors to investigate the rank of a matrix. In particular, we have the following proposition.

**Proposition 2.4** *If an  $m \times n$  ( $m \geq n$ ) matrix  $\mathbf{A}$  has a nonzero  $n$ th-order minor, then the columns of  $\mathbf{A}$  are linearly independent; that is,  $\text{rank } \mathbf{A} = n$ .*

□

*Proof.* Suppose that  $\mathbf{A}$  has a nonzero  $n$ th-order minor. Without loss of generality, we assume that the  $n$ th-order minor corresponding to the first  $n$  rows of  $\mathbf{A}$  is nonzero. Let  $x_i$ ,  $i = 1, \dots, n$ , be scalars such that

$$x_1 a_1 + x_2 a_2 + \cdots + x_n a_n = \mathbf{0}.$$

The vector equality above is equivalent to the following set of  $m$  equations:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= 0 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= 0 \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= 0 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= 0. \end{aligned}$$

For  $i = 1, \dots, n$ , let

$$\tilde{\mathbf{a}}_i = \begin{bmatrix} a_{1i} \\ \vdots \\ a_{ni} \end{bmatrix}.$$

Then,  $x_1\tilde{\mathbf{a}}_1 + \cdots + x_n\tilde{\mathbf{a}}_n = \mathbf{0}$ .

The  $n$ th-order minor is  $\det[\tilde{\mathbf{a}}_1, \tilde{\mathbf{a}}_2, \dots, \tilde{\mathbf{a}}_n]$ , assumed to be nonzero. From the properties of determinants it follows that the columns  $\tilde{\mathbf{a}}_1, \tilde{\mathbf{a}}_2, \dots, \tilde{\mathbf{a}}_n$  are linearly independent. Therefore, all  $x_i = 0$ ,  $i = 1, \dots, n$ . Hence, the columns  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$  are linearly independent. ■

From the above it follows that if there is a nonzero minor, then the columns associated with this nonzero minor are linearly independent.

If a matrix  $\mathbf{A}$  has an  $r$ th-order minor  $|\mathbf{M}|$  with the properties (i)  $|\mathbf{M}| \neq 0$  and (ii) any minor of  $\mathbf{A}$  that is formed by adding a row and a column of  $\mathbf{A}$  to  $\mathbf{M}$  is zero, then

$$\text{rank } \mathbf{A} = r.$$

Thus, the rank of a matrix is equal to the highest order of its nonzero minor(s).

A *nonsingular* (or *invertible*) matrix is a square matrix whose determinant is nonzero. Suppose that  $\mathbf{A}$  is an  $n \times n$  square matrix. Then,  $\mathbf{A}$  is nonsingular if and only if there is another  $n \times n$  matrix  $\mathbf{B}$  such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I}_n,$$

where  $\mathbf{I}_n$  denotes the  $n \times n$  *identity matrix*:

$$\mathbf{I}_n = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}.$$

We call the matrix  $\mathbf{B}$  above the *inverse matrix* of  $\mathbf{A}$ , and write  $\mathbf{B} = \mathbf{A}^{-1}$ .

### 2.3 Linear Equations

Suppose that we are given  $m$  equations in  $n$  unknowns of the form

$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1,$$

$$a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2,$$

 $\vdots$ 

$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m.$$

We can represent the set of equations above as a vector equation

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b},$$

where

$$\mathbf{a}_j = \begin{bmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}.$$

Associated with this system of equations is the matrix

$$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n],$$

and an augmented matrix

$$[\mathbf{A}, \mathbf{b}] = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n, \mathbf{b}].$$

We can also represent the system of equations above as

$$\mathbf{Ax} = \mathbf{b},$$

where

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

**Theorem 2.1** *The system of equations  $\mathbf{Ax} = \mathbf{b}$  has a solution if and only if*

$$\text{rank } \mathbf{A} = \text{rank}[\mathbf{A}, \mathbf{b}].$$

□

*Proof.*  $\Rightarrow$ : Suppose that the system  $\mathbf{Ax} = \mathbf{b}$  has a solution. Therefore,  $\mathbf{b}$  is a linear combination of the columns of  $\mathbf{A}$ ; that is, there exist  $x_1, \dots, x_n$  such

that  $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b}$ . It follows that  $\mathbf{b}$  belongs to  $\text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n]$  and hence

$$\begin{aligned}\text{rank } \mathbf{A} &= \dim \text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n] \\ &= \dim \text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{b}] \\ &= \text{rank}[\mathbf{A}, \mathbf{b}].\end{aligned}$$

$\Leftarrow$ : Suppose that  $\text{rank } \mathbf{A} = \text{rank}[\mathbf{A}, \mathbf{b}] = r$ . Thus, we have  $r$  linearly independent columns of  $\mathbf{A}$ . Without loss of generality, let  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$  be these columns. Therefore,  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$  are also linearly independent columns of the matrix  $[\mathbf{A}, \mathbf{b}]$ . Because  $\text{rank}[\mathbf{A}, \mathbf{b}] = r$ , the remaining columns of  $[\mathbf{A}, \mathbf{b}]$  can be expressed as linear combinations of  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$ . In particular,  $\mathbf{b}$  can be expressed as a linear combination of these columns. Hence, there exist  $x_1, \dots, x_n$  such that  $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_n\mathbf{a}_n = \mathbf{b}$ . ■

**Theorem 2.2** Consider the equation  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\text{rank } \mathbf{A} = m$ . A solution to  $\mathbf{Ax} = \mathbf{b}$  can be obtained by assigning arbitrary values for  $n - m$  variables and solving for the remaining ones. □

*Proof.* We have  $\text{rank } \mathbf{A} = m$ , and therefore we can find  $m$  linearly independent columns of  $\mathbf{A}$ . Without loss of generality, let  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m$  be such columns. Rewrite the equation  $\mathbf{Ax} = \mathbf{b}$  as

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_m\mathbf{a}_m = \mathbf{b} - x_{m+1}\mathbf{a}_{m+1} - \cdots - x_n\mathbf{a}_n.$$

Assign to  $x_{m+1}, x_{m+2}, \dots, x_n$  arbitrary values, say

$$x_{m+1} = d_{m+1}, x_{m+2} = d_{m+2}, \dots, x_n = d_n,$$

and let

$$\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m] \in \mathbb{R}^{m \times m}.$$

Note that  $\det \mathbf{B} \neq 0$ . We can represent the system of equations above as

$$\mathbf{B} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = [\mathbf{b} - d_{m+1}\mathbf{a}_{m+1} - \cdots - d_n\mathbf{a}_n].$$

The matrix  $\mathbf{B}$  is invertible, and therefore we can solve for  $[x_1, x_2, \dots, x_m]^\top$ . Specifically,

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \mathbf{B}^{-1} [\mathbf{b} - d_{m+1}\mathbf{a}_{m+1} - \cdots - d_n\mathbf{a}_n].$$

■

## 2.4 Inner Products and Norms

The absolute value of a real number  $a$ , denoted  $|a|$ , is defined as

$$|a| = \begin{cases} a & \text{if } a \geq 0 \\ -a & \text{if } a < 0. \end{cases}$$

The following formulas hold:

1.  $|a| = |-a|$ .
2.  $-|a| \leq a \leq |a|$ .
3.  $|a + b| \leq |a| + |b|$ .
4.  $||a| - |b|| \leq |a - b| \leq |a| + |b|$ .
5.  $|ab| = |a||b|$ .
6.  $|a| \leq c$  and  $|b| \leq d$  imply that  $|a + b| \leq c + d$ .
7. The inequality  $|a| < b$  is equivalent to  $-b < a < b$  (i.e.,  $a < b$  and  $-a < b$ ). The same holds if we replace every occurrence of " $<$ " by " $\leq$ ".
8. The inequality  $|a| > b$  is equivalent to  $a > b$  or  $-a > b$ . The same holds if we replace every occurrence of " $>$ " by " $\geq$ ".

For  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , we define the *Euclidean inner product* by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n x_i y_i = \mathbf{x}^\top \mathbf{y}.$$

The inner product is a real-valued function  $\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  having the following properties:

1. Positivity:  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ ,  $\langle \mathbf{x}, \mathbf{x} \rangle = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ .
2. Symmetry:  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ .
3. Additivity:  $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$ .
4. Homogeneity:  $\langle r\mathbf{x}, \mathbf{y} \rangle = r\langle \mathbf{x}, \mathbf{y} \rangle$  for every  $r \in \mathbb{R}$ .

The properties of additivity and homogeneity in the second vector also hold; that is,

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle &= \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle, \\ \langle \mathbf{x}, r\mathbf{y} \rangle &= r\langle \mathbf{x}, \mathbf{y} \rangle \quad \text{for every } r \in \mathbb{R}. \end{aligned}$$

The above can be shown using properties 2 to 4. Indeed,

$$\begin{aligned}\langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle &= \langle \mathbf{y} + \mathbf{z}, \mathbf{x} \rangle \\ &= \langle \mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{z}, \mathbf{x} \rangle \\ &= \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle\end{aligned}$$

and

$$\langle \mathbf{x}, r\mathbf{y} \rangle = \langle r\mathbf{y}, \mathbf{x} \rangle = r\langle \mathbf{y}, \mathbf{x} \rangle = r\langle \mathbf{x}, \mathbf{y} \rangle.$$

It is possible to define other real-valued functions on  $\mathbb{R}^n \times \mathbb{R}^n$  that satisfy properties 1 to 4 above (see Exercise 2.8). Many results involving the Euclidean inner product also hold for these other forms of inner products.

The vectors  $\mathbf{x}$  and  $\mathbf{y}$  are said to be *orthogonal* if  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ .

The *Euclidean norm* of a vector  $\mathbf{x}$  is defined as

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\mathbf{x}^\top \mathbf{x}}.$$

**Theorem 2.3 Cauchy-Schwarz Inequality.** *For any two vectors  $\mathbf{x}$  and  $\mathbf{y}$  in  $\mathbb{R}^n$ , the Cauchy-Schwarz inequality*

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

*holds. Furthermore, equality holds if and only if  $\mathbf{x} = \alpha\mathbf{y}$  for some  $\alpha \in \mathbb{R}$ .*  $\square$

*Proof.* First assume that  $\mathbf{x}$  and  $\mathbf{y}$  are unit vectors; that is,  $\|\mathbf{x}\| = \|\mathbf{y}\| = 1$ . Then,

$$\begin{aligned}0 \leq \|\mathbf{x} - \mathbf{y}\|^2 &= \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle \\ &= \|\mathbf{x}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 \\ &= 2 - 2\langle \mathbf{x}, \mathbf{y} \rangle\end{aligned}$$

or

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq 1,$$

with equality holding if and only if  $\mathbf{x} = \mathbf{y}$ .

Next, assuming that neither  $\mathbf{x}$  nor  $\mathbf{y}$  is zero (for the inequality obviously holds if one of them is zero), we replace  $\mathbf{x}$  and  $\mathbf{y}$  by the unit vectors  $\mathbf{x}/\|\mathbf{x}\|$  and  $\mathbf{y}/\|\mathbf{y}\|$ . Then, apply property 4 to get

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

Now replace  $\mathbf{x}$  by  $-\mathbf{x}$  and again apply property 4 to get

$$-\langle \mathbf{x}, \mathbf{y} \rangle \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

The last two inequalities imply the absolute value inequality. Equality holds if and only if  $\mathbf{x}/\|\mathbf{x}\| = \pm \mathbf{y}/\|\mathbf{y}\|$ ; that is,  $\mathbf{x} = \alpha\mathbf{y}$  for some  $\alpha \in \mathbb{R}$ .  $\blacksquare$

The Euclidean norm of a vector  $\|\mathbf{x}\|$  has the following properties:

1. Positivity:  $\|\mathbf{x}\| \geq 0$ ,  $\|\mathbf{x}\| = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ .
2. Homogeneity:  $\|r\mathbf{x}\| = |r|\|\mathbf{x}\|$ ,  $r \in \mathbb{R}$ .
3. Triangle inequality:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ .

The triangle inequality can be proved using the Cauchy-Schwarz inequality, as follows. We have

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2.$$

By the Cauchy-Schwarz inequality,

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^2 &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 \\ &= (\|\mathbf{x}\| + \|\mathbf{y}\|)^2, \end{aligned}$$

and therefore

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|.$$

Note that if  $\mathbf{x}$  and  $\mathbf{y}$  are orthogonal:  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ , then

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2,$$

which is the *Pythagorean theorem* for  $\mathbb{R}^n$ .

The Euclidean norm is an example of a general *vector norm*, which is any function satisfying the three properties of positivity, homogeneity, and triangle inequality. Other examples of vector norms on  $\mathbb{R}^n$  include the 1-norm, defined by  $\|\mathbf{x}\|_1 = |x_1| + \cdots + |x_n|$ , and the  $\infty$ -norm, defined by  $\|\mathbf{x}\|_\infty = \max_i |x_i|$  (where the notation  $\max_i$  represents the largest over all the possible index values of  $i$ ). The Euclidean norm is often referred to as the *2-norm*, and denoted  $\|\mathbf{x}\|_2$ . The norms above are special cases of the *p-norm*, given by

$$\|\mathbf{x}\|_p = \begin{cases} (|x_1|^p + \cdots + |x_n|^p)^{1/p} & \text{if } 1 \leq p < \infty \\ \max\{|x_1|, \dots, |x_n|\} & \text{if } p = \infty. \end{cases}$$

We can use norms to define the notion of a continuous function, as follows. A function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is *continuous* at  $\mathbf{x}$  if for all  $\varepsilon > 0$ , there exists  $\delta > 0$  such that  $\|\mathbf{y} - \mathbf{x}\| < \delta \Rightarrow \|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})\| < \varepsilon$ . If the function  $\mathbf{f}$  is continuous at every point in  $\mathbb{R}^n$ , we say that it is continuous on  $\mathbb{R}^n$ . Note that  $\mathbf{f} = [f_1, \dots, f_m]^\top$  is continuous if and only if each component  $f_i$ ,  $i = 1, \dots, m$ , is continuous.

For the complex vector space  $\mathbb{C}^n$ , we define an inner product  $\langle \mathbf{x}, \mathbf{y} \rangle$  to be  $\sum_{i=1}^n x_i \bar{y}_i$ , where the bar denotes complex conjugation. The inner product on  $\mathbb{C}^n$  is a complex-valued function having the following properties:

1.  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ ,  $\langle \mathbf{x}, \mathbf{x} \rangle = 0$  if and only if  $\mathbf{x} = \mathbf{0}$ .

2.  $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}.$
3.  $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle.$
4.  $\langle r\mathbf{x}, \mathbf{y} \rangle = r\langle \mathbf{x}, \mathbf{y} \rangle$ , where  $r \in \mathbb{C}.$

From properties 1 to 4, we can deduce other properties, such as

$$\langle \mathbf{x}, r_1\mathbf{y} + r_2\mathbf{z} \rangle = \bar{r}_1\langle \mathbf{x}, \mathbf{y} \rangle + \bar{r}_2\langle \mathbf{x}, \mathbf{z} \rangle,$$

where  $r_1, r_2 \in \mathbb{C}$ . For  $\mathbb{C}^n$ , the vector norm can similarly be defined by  $\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle$ . For more information, consult Gel'fand [47].

## EXERCISES

**2.1** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\text{rank } \mathbf{A} = m$ . Show that  $m \leq n$ .

**2.2** Prove that the system  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , has a unique solution if and only if  $\text{rank } \mathbf{A} = \text{rank}[\mathbf{A}, \mathbf{b}] = n$ .

**2.3** (Adapted from [38].) We know that if  $k \geq n + 1$ , then the vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k \in \mathbb{R}^n$  are linearly dependent; that is, there exist scalars  $\alpha_1, \dots, \alpha_k$  such that at least one  $\alpha_i \neq 0$  and  $\sum_{i=1}^k \alpha_i \mathbf{a}_i = \mathbf{0}$ . Show that if  $k \geq n + 2$ , then there exist scalars  $\alpha_1, \dots, \alpha_k$  such that at least one  $\alpha_i \neq 0$ ,  $\sum_{i=1}^k \alpha_i \mathbf{a}_i = \mathbf{0}$ , and  $\sum_{i=1}^k \alpha_i = 0$ .

*Hint:* Introduce the vectors  $\bar{\mathbf{a}}_i = [1, \mathbf{a}_i^\top]^\top \in \mathbb{R}^{n+1}$ ,  $i = 1, \dots, k$ , and use the fact that any  $n + 2$  vectors in  $\mathbb{R}^{n+1}$  are linearly dependent.

**2.4** Consider an  $m \times m$  matrix  $\mathbf{M}$  that has block form

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{m-k,k} & \mathbf{I}_{m-k} \\ \mathbf{M}_{k,k} & \mathbf{O}_{k,m-k} \end{bmatrix},$$

where  $\mathbf{M}_{k,k}$  is  $k \times k$ ,  $\mathbf{M}_{m-k,k}$  is  $(m-k) \times k$ ,  $\mathbf{I}_{m-k}$  is the  $(m-k) \times (m-k)$  identity matrix, and  $\mathbf{O}_{k,m-k}$  is the  $k \times (m-k)$  zero matrix.

- a. Show that

$$|\det \mathbf{M}| = |\det \mathbf{M}_{k,k}|.$$

This result is relevant to the proof of Proposition 19.1.

- b. Under certain assumptions, the following stronger result holds:

$$\det \mathbf{M} = \det(-\mathbf{M}_{k,k})$$

Identify cases where this is true, and show that it is false in general.

**2.5** It is well known that for any  $a, b, c, d \in \mathbb{C}$ ,

$$\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc.$$

Suppose now that  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{D}$  are real or complex square matrices of the same size. Give a sufficient condition under which

$$\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \mathbf{AD} - \mathbf{BC}.$$

An interesting discussion on determinants of block matrices is provided in [121].

**2.6** Consider the following system of linear equations:

$$\begin{aligned} x_1 + x_2 + 2x_3 + x_4 &= 1 \\ x_1 - 2x_2 - x_4 &= -2. \end{aligned}$$

Use Theorem 2.1 to check if the system has a solution. Then, use the method of Theorem 2.2 to find a general solution to the system.

**2.7** Prove the seven properties of the absolute value of a real number.

**2.8** Consider the function  $\langle \cdot, \cdot \rangle_2 : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ , defined by  $\langle \mathbf{x}, \mathbf{y} \rangle_2 = 2x_1y_1 + 3x_2y_1 + 3x_1y_2 + 5x_2y_2$ , where  $\mathbf{x} = [x_1, x_2]^\top$  and  $\mathbf{y} = [y_1, y_2]^\top$ . Show that  $\langle \cdot, \cdot \rangle_2$  satisfies conditions 1 to 4 for inner products.

*Note:* This is a special case of Exercise 3.21.

**2.9** Show that for any two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ,  $|||\mathbf{x}| - |\mathbf{y}||| \leq \|\mathbf{x} - \mathbf{y}\|$ .

*Hint:* Write  $\mathbf{x} = (\mathbf{x} - \mathbf{y}) + \mathbf{y}$ , and use the triangle inequality. Do the same for  $\mathbf{y}$ .

**2.10** Use Exercise 2.9 to show that the norm  $\|\cdot\|$  is a *uniformly continuous function*; that is, for all  $\varepsilon > 0$ , there exists  $\delta > 0$  such that if  $\|\mathbf{x} - \mathbf{y}\| < \delta$ , then  $|||\mathbf{x}| - |\mathbf{y}||| < \varepsilon$ .



# CHAPTER 3

---

## TRANSFORMATIONS

---

### 3.1 Linear Transformations

A function  $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is called a *linear transformation* if:

1.  $\mathcal{L}(ax) = a\mathcal{L}(x)$  for every  $x \in \mathbb{R}^n$  and  $a \in \mathbb{R}$ .
2.  $\mathcal{L}(x + y) = \mathcal{L}(x) + \mathcal{L}(y)$  for every  $x, y \in \mathbb{R}^n$ .

If we fix the bases for  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , then the linear transformation  $\mathcal{L}$  can be represented by a matrix. Specifically, there exists  $\mathbf{A} \in \mathbb{R}^{m \times n}$  such that the following representation holds. Suppose that  $\mathbf{x} \in \mathbb{R}^n$  is a given vector, and  $\mathbf{x}'$  is the representation of  $\mathbf{x}$  with respect to the given basis for  $\mathbb{R}^n$ . If  $\mathbf{y} = \mathcal{L}(\mathbf{x})$ , and  $\mathbf{y}'$  is the representation of  $\mathbf{y}$  with respect to the given basis for  $\mathbb{R}^m$ , then

$$\mathbf{y}' = \mathbf{Ax}'.$$

We call  $\mathbf{A}$  the *matrix representation* of  $\mathcal{L}$  with respect to the given bases for  $\mathbb{R}^n$  and  $\mathbb{R}^m$ . In the special case where we assume the natural bases for  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , the matrix representation  $\mathbf{A}$  satisfies

$$\mathcal{L}(\mathbf{x}) = \mathbf{Ax}.$$

Let  $\{e_1, e_2, \dots, e_n\}$  and  $\{e'_1, e'_2, \dots, e'_n\}$  be two bases for  $\mathbb{R}^n$ . Define the matrix

$$T = [e'_1, e'_2, \dots, e'_n]^{-1}[e_1, e_2, \dots, e_n].$$

We call  $T$  the *transformation matrix* from  $\{e_1, e_2, \dots, e_n\}$  to  $\{e'_1, e'_2, \dots, e'_n\}$ . It is clear that

$$[e_1, e_2, \dots, e_n] = [e'_1, e'_2, \dots, e'_n]T;$$

that is, the  $i$ th column of  $T$  is the vector of coordinates of  $e_i$  with respect to the basis  $\{e'_1, e'_2, \dots, e'_n\}$ .

Fix a vector in  $\mathbb{R}^n$ . Let  $x$  be the column of the coordinates of the vector with respect to  $\{e_1, e_2, \dots, e_n\}$  and  $x'$  the coordinates of the same vector with respect to  $\{e'_1, e'_2, \dots, e'_n\}$ . Then, we can show that  $x' = Tx$  (see Exercise 3.1).

Consider a linear transformation

$$\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}^n,$$

and let  $A$  be its representation with respect to  $\{e_1, e_2, \dots, e_n\}$  and  $B$  its representation with respect to  $\{e'_1, e'_2, \dots, e'_n\}$ . Let  $y = Ax$  and  $y' = Bx'$ . Therefore,  $y' = Ty = TAx = Bx' = BTx$ , and hence  $TA = BT$ , or  $A = T^{-1}BT$ .

Two  $n \times n$  matrices  $A$  and  $B$  are *similar* if there exists a nonsingular matrix  $T$  such that  $A = T^{-1}BT$ . In conclusion, similar matrices correspond to the same linear transformation with respect to different bases.

### 3.2 Eigenvalues and Eigenvectors

Let  $A$  be an  $n \times n$  real square matrix. A scalar  $\lambda$  (possibly complex) and a nonzero vector  $v$  satisfying the equation  $Av = \lambda v$  are said to be, respectively, an *eigenvalue* and an *eigenvector* of  $A$ . For  $\lambda$  to be an eigenvalue it is necessary and sufficient for the matrix  $\lambda I - A$  to be singular; that is,  $\det[\lambda I - A] = 0$ , where  $I$  is the  $n \times n$  identity matrix. This leads to an  $n$ th-order polynomial equation

$$\det[\lambda I - A] = \lambda^n + a_{n-1}\lambda^{n-1} + \cdots + a_1\lambda + a_0 = 0.$$

We call the polynomial  $\det[\lambda I - A]$  the *characteristic polynomial* of the matrix  $A$ , and the equation above the *characteristic equation*. According to the fundamental theorem of algebra, the characteristic equation must have  $n$  (possibly nondistinct) roots that are the eigenvalues of  $A$ . The following theorem states that if  $A$  has  $n$  distinct eigenvalues, then it also has  $n$  linearly independent eigenvectors.

**Theorem 3.1** Suppose that the characteristic equation  $\det[\lambda I - A] = 0$  has  $n$  distinct roots  $\lambda_1, \lambda_2, \dots, \lambda_n$ . Then, there exist  $n$  linearly independent vectors  $v_1, v_2, \dots, v_n$  such that

$$Av_i = \lambda_i v_i, \quad i = 1, 2, \dots, n.$$

□

*Proof.* Because  $\det[\lambda_i \mathbf{I} - \mathbf{A}] = 0$ ,  $i = 1, \dots, n$ , there exist nonzero  $\mathbf{v}_i$ ,  $i = 1, \dots, n$ , such that  $\mathbf{A}\mathbf{v}_i = \lambda_i \mathbf{v}_i$ ,  $i = 1, \dots, n$ . We now prove the linear independence of  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ . To do this, let  $c_1, \dots, c_n$  be scalars such that  $\sum_{i=1}^n c_i \mathbf{v}_i = \mathbf{0}$ . We show that  $c_i = 0$ ,  $i = 1, \dots, n$ .

Consider the matrix

$$\mathbf{Z} = (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A}) \cdots (\lambda_n \mathbf{I} - \mathbf{A}).$$

We first show that  $c_1 = 0$ . Note that

$$\begin{aligned} \mathbf{Z}\mathbf{v}_n &= (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A}) \cdots (\lambda_{n-1} \mathbf{I} - \mathbf{A})(\lambda_n \mathbf{I} - \mathbf{A})\mathbf{v}_n \\ &= (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A}) \cdots (\lambda_{n-1} \mathbf{I} - \mathbf{A})(\lambda_n \mathbf{v}_n - \mathbf{A}\mathbf{v}_n) \\ &= \mathbf{0} \end{aligned}$$

since  $\lambda_n \mathbf{v}_n - \mathbf{A}\mathbf{v}_n = \mathbf{0}$ .

Repeating the argument above, we get

$$\mathbf{Z}\mathbf{v}_k = \mathbf{0}, \quad k = 2, 3, \dots, n.$$

But

$$\begin{aligned} \mathbf{Z}\mathbf{v}_1 &= (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A}) \cdots (\lambda_{n-1} \mathbf{I} - \mathbf{A})(\lambda_n \mathbf{I} - \mathbf{A})\mathbf{v}_1 \\ &= (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A}) \cdots (\lambda_{n-1} \mathbf{v}_1 - \mathbf{A}\mathbf{v}_1)(\lambda_n - \lambda_1) \\ &\quad \vdots \\ &= (\lambda_2 \mathbf{I} - \mathbf{A})(\lambda_3 \mathbf{I} - \mathbf{A})\mathbf{v}_1 \cdots (\lambda_{n-1} - \lambda_1)(\lambda_n - \lambda_1) \\ &= (\lambda_2 - \lambda_1)(\lambda_3 - \lambda_1) \cdots (\lambda_{n-1} - \lambda_1)(\lambda_n - \lambda_1)\mathbf{v}_1. \end{aligned}$$

Using the equation above, we see that

$$\begin{aligned} \mathbf{Z} \left( \sum_{i=1}^n c_i \mathbf{v}_i \right) &= \sum_{i=1}^n c_i \mathbf{Z}\mathbf{v}_i \\ &= c_1 \mathbf{Z}\mathbf{v}_1 \\ &= c_1 (\lambda_2 - \lambda_1)(\lambda_3 - \lambda_1) \cdots (\lambda_n - \lambda_1)\mathbf{v}_1 = \mathbf{0}. \end{aligned}$$

Because the  $\lambda_i$  are distinct, it must follow that  $c_1 = 0$ .

Using similar arguments, we can show that all  $c_i$  must vanish, and therefore the set of eigenvectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  is linearly independent. ■

Consider a basis formed by a linearly independent set of eigenvectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ . With respect to this basis, the matrix  $\mathbf{A}$  is *diagonal* [i.e., if  $a_{ij}$  is the  $(i, j)$ th element of  $\mathbf{A}$ , then  $a_{ij} = 0$  for all  $i \neq j$ ]. Indeed, let

$$\mathbf{T} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]^{-1}.$$

Then,

$$\begin{aligned}
 \mathbf{T}A\mathbf{T}^{-1} &= \mathbf{T}\mathbf{A}[v_1, v_2, \dots, v_n] \\
 &= \mathbf{T}[\mathbf{A}v_1, \mathbf{A}v_2, \dots, \mathbf{A}v_n] \\
 &= \mathbf{T}[\lambda_1 v_1, \lambda_2 v_2, \dots, \lambda_n v_n] \\
 &= \mathbf{T}\mathbf{T}^{-1} \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_n \end{bmatrix} \\
 &= \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_n \end{bmatrix},
 \end{aligned}$$

because  $\mathbf{T}\mathbf{T}^{-1} = \mathbf{I}$ .

A matrix  $\mathbf{A}$  is *symmetric* if  $\mathbf{A} = \mathbf{A}^\top$ .

**Theorem 3.2** All eigenvalues of a real symmetric matrix are real. □

*Proof.* Let

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x},$$

where  $\mathbf{x} \neq \mathbf{0}$ . Taking the inner product of  $\mathbf{A}\mathbf{x}$  with  $\mathbf{x}$  yields

$$\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle = \langle \lambda\mathbf{x}, \mathbf{x} \rangle = \lambda\langle \mathbf{x}, \mathbf{x} \rangle.$$

On the other hand,

$$\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{A}^\top \mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \langle \mathbf{x}, \lambda\mathbf{x} \rangle = \bar{\lambda}\langle \mathbf{x}, \mathbf{x} \rangle.$$

The above follows from the definition of the inner product on  $\mathbb{C}^n$ . We note that  $\langle \mathbf{x}, \mathbf{x} \rangle$  is real and  $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ . Hence,

$$\lambda\langle \mathbf{x}, \mathbf{x} \rangle = \bar{\lambda}\langle \mathbf{x}, \mathbf{x} \rangle$$

and

$$(\lambda - \bar{\lambda})\langle \mathbf{x}, \mathbf{x} \rangle = 0.$$

Because  $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ ,

$$\lambda = \bar{\lambda}.$$

Thus,  $\lambda$  is real. ■

**Theorem 3.3** Any real symmetric  $n \times n$  matrix has a set of  $n$  eigenvectors that are mutually orthogonal. □

*Proof.* We prove the result for the case when the  $n$  eigenvalues are distinct. For a general proof, see [62, p. 104].

Suppose that  $\mathbf{A}\mathbf{v}_1 = \lambda_1\mathbf{v}_1$ ,  $\mathbf{A}\mathbf{v}_2 = \lambda_2\mathbf{v}_2$ , where  $\lambda_1 \neq \lambda_2$ . Then,

$$\langle \mathbf{A}\mathbf{v}_1, \mathbf{v}_2 \rangle = \langle \lambda_1\mathbf{v}_1, \mathbf{v}_2 \rangle = \lambda_1 \langle \mathbf{v}_1, \mathbf{v}_2 \rangle.$$

Because  $\mathbf{A} = \mathbf{A}^\top$ ,

$$\langle \mathbf{A}\mathbf{v}_1, \mathbf{v}_2 \rangle = \langle \mathbf{v}_1, \mathbf{A}^\top \mathbf{v}_2 \rangle = \langle \mathbf{v}_1, \mathbf{A}\mathbf{v}_2 \rangle = \lambda_2 \langle \mathbf{v}_1, \mathbf{v}_2 \rangle.$$

Therefore,

$$\lambda_1 \langle \mathbf{v}_1, \mathbf{v}_2 \rangle = \lambda_2 \langle \mathbf{v}_1, \mathbf{v}_2 \rangle.$$

Because  $\lambda_1 \neq \lambda_2$ , it follows that

$$\langle \mathbf{v}_1, \mathbf{v}_2 \rangle = 0.$$

■

If  $\mathbf{A}$  is symmetric, then a set of its eigenvectors forms an orthogonal basis for  $\mathbb{R}^n$ . If the basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  is normalized so that each element has norm of unity, then defining the matrix

$$\mathbf{T} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n],$$

we have

$$\mathbf{T}^\top \mathbf{T} = \mathbf{I}$$

and hence

$$\mathbf{T}^\top = \mathbf{T}^{-1}.$$

A matrix whose transpose is its inverse is said to be an *orthogonal matrix*.

### 3.3 Orthogonal Projections

Recall that a subspace  $\mathcal{V}$  of  $\mathbb{R}^n$  is a subset that is closed under the operations of vector addition and scalar multiplication. In other words,  $\mathcal{V}$  is a subspace of  $\mathbb{R}^n$  if  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{V} \Rightarrow \alpha\mathbf{x}_1 + \beta\mathbf{x}_2 \in \mathcal{V}$  for all  $\alpha, \beta \in \mathbb{R}$ . Furthermore, the dimension of a subspace  $\mathcal{V}$  is equal to the maximum number of linearly independent vectors in  $\mathcal{V}$ . If  $\mathcal{V}$  is a subspace of  $\mathbb{R}^n$ , then the *orthogonal complement* of  $\mathcal{V}$ , denoted  $\mathcal{V}^\perp$ , consists of all vectors that are orthogonal to every vector in  $\mathcal{V}$ . Thus,

$$\mathcal{V}^\perp = \{\mathbf{x} : \mathbf{v}^\top \mathbf{x} = 0 \text{ for all } \mathbf{v} \in \mathcal{V}\}.$$

The orthogonal complement of  $\mathcal{V}$  is also a subspace (see Exercise 3.7). Together,  $\mathcal{V}$  and  $\mathcal{V}^\perp$  span  $\mathbb{R}^n$  in the sense that every vector  $\mathbf{x} \in \mathbb{R}^n$  can be represented uniquely as

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2,$$

where  $\mathbf{x}_1 \in \mathcal{V}$  and  $\mathbf{x}_2 \in \mathcal{V}^\perp$ . We call the representation above the *orthogonal decomposition* of  $\mathbf{x}$  (with respect to  $\mathcal{V}$ ). We say that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are *orthogonal projections* of  $\mathbf{x}$  onto the subspaces  $\mathcal{V}$  and  $\mathcal{V}^\perp$ , respectively. We write  $\mathbb{R}^n = \mathcal{V} \oplus \mathcal{V}^\perp$  and say that  $\mathbb{R}^n$  is a *direct sum* of  $\mathcal{V}$  and  $\mathcal{V}^\perp$ . We say that a linear transformation  $\mathbf{P}$  is an *orthogonal projector* onto  $\mathcal{V}$  if for all  $\mathbf{x} \in \mathbb{R}^n$ , we have  $\mathbf{P}\mathbf{x} \in \mathcal{V}$  and  $\mathbf{x} - \mathbf{P}\mathbf{x} \in \mathcal{V}^\perp$ .

In the subsequent discussion we use the following notation. Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Let the *range*, or *image*, of  $\mathbf{A}$  be denoted

$$\mathcal{R}(\mathbf{A}) \triangleq \{\mathbf{Ax} : \mathbf{x} \in \mathbb{R}^n\},$$

and the *nullspace*, or *kernel*, of  $\mathbf{A}$  be denoted

$$\mathcal{N}(\mathbf{A}) \triangleq \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{0}\}.$$

Note that  $\mathcal{R}(\mathbf{A})$  and  $\mathcal{N}(\mathbf{A})$  are subspaces (see Exercise 3.9).

**Theorem 3.4** *Let  $\mathbf{A}$  be a given matrix. Then,  $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^\top)$  and  $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top)$ .*  $\square$

*Proof.* Suppose that  $\mathbf{x} \in \mathcal{R}(\mathbf{A})^\perp$ . Then,  $\mathbf{y}^\top(\mathbf{A}^\top \mathbf{x}) = (\mathbf{Ay})^\top \mathbf{x} = 0$  for all  $\mathbf{y}$ , so that  $\mathbf{A}^\top \mathbf{x} = \mathbf{0}$ . Hence,  $\mathbf{x} \in \mathcal{N}(\mathbf{A}^\top)$ . This implies that  $\mathcal{R}(\mathbf{A})^\perp \subset \mathcal{N}(\mathbf{A}^\top)$ .

If now  $\mathbf{x} \in \mathcal{N}(\mathbf{A}^\top)$ , then  $(\mathbf{Ay})^\top \mathbf{x} = \mathbf{y}^\top(\mathbf{A}^\top \mathbf{x}) = 0$  for all  $\mathbf{y}$ , so that  $\mathbf{x} \in \mathcal{R}(\mathbf{A})^\perp$ , and consequently,  $\mathcal{N}(\mathbf{A}^\top) \subset \mathcal{R}(\mathbf{A})^\perp$ . Thus,  $\mathcal{R}(\mathbf{A})^\perp = \mathcal{N}(\mathbf{A}^\top)$ .

The equation  $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top)$  follows from what we have proved above and the fact that for any subspace  $\mathcal{V}$ , we have  $(\mathcal{V}^\perp)^\perp = \mathcal{V}$  (see Exercise 3.11).  $\blacksquare$

Theorem 3.4 allows us to establish the following necessary and sufficient condition for orthogonal projectors. For this, note that if  $\mathbf{P}$  is an orthogonal projector onto  $\mathcal{V}$ , then  $\mathbf{Px} = \mathbf{x}$  for all  $\mathbf{x} \in \mathcal{V}$ , and  $\mathcal{R}(\mathbf{P}) = \mathcal{V}$  (see Exercise 3.14).

**Theorem 3.5** *A matrix  $\mathbf{P}$  is an orthogonal projector [onto the subspace  $\mathcal{V} = \mathcal{R}(\mathbf{P})$ ] if and only if  $\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^\top$ .*  $\square$

*Proof.*  $\Rightarrow$ : Suppose that  $\mathbf{P}$  is an orthogonal projector onto  $\mathcal{V} = \mathcal{R}(\mathbf{P})$ . Then,  $\mathcal{R}(\mathbf{I} - \mathbf{P}) \subset \mathcal{R}(\mathbf{P})^\perp$ . But,  $\mathcal{R}(\mathbf{P})^\perp = \mathcal{N}(\mathbf{P}^\top)$  by Theorem 3.4. Therefore,  $\mathcal{R}(\mathbf{I} - \mathbf{P}) \subset \mathcal{N}(\mathbf{P}^\top)$ . Hence,  $\mathbf{P}^\top(\mathbf{I} - \mathbf{P})\mathbf{y} = \mathbf{0}$  for all  $\mathbf{y}$ , which implies that  $\mathbf{P}^\top(\mathbf{I} - \mathbf{P}) = \mathbf{O}$ , where  $\mathbf{O}$  is the matrix with all entries equal to zero; i.e., the *zero matrix*. Therefore,  $\mathbf{P}^\top = \mathbf{P}^\top \mathbf{P}$ , and thus  $\mathbf{P} = \mathbf{P}^\top = \mathbf{P}^2$ .

$\Leftarrow$ : Suppose that  $\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^\top$ . For any  $\mathbf{x}$ , we have  $(\mathbf{Py})^\top(\mathbf{I} - \mathbf{P})\mathbf{x} = \mathbf{y}^\top \mathbf{P}^\top(\mathbf{I} - \mathbf{P})\mathbf{x} = \mathbf{y}^\top \mathbf{P}(\mathbf{I} - \mathbf{P})\mathbf{x} = 0$  for all  $\mathbf{y}$ . Thus,  $(\mathbf{I} - \mathbf{P})\mathbf{x} \in \mathcal{R}(\mathbf{P})^\perp$ , which means that  $\mathbf{P}$  is an orthogonal projector.  $\blacksquare$

### 3.4 Quadratic Forms

A *quadratic form*  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a function

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x},$$

where  $\mathbf{Q}$  is an  $n \times n$  real matrix. There is no loss of generality in assuming  $\mathbf{Q}$  to be symmetric:  $\mathbf{Q} = \mathbf{Q}^\top$ . For if the matrix  $\mathbf{Q}$  is not symmetric, we can always replace it with the symmetric matrix

$$\mathbf{Q}_0 = \mathbf{Q}_0^\top = \frac{1}{2} (\mathbf{Q} + \mathbf{Q}^\top).$$

Note that

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \mathbf{x}^\top \mathbf{Q}_0 \mathbf{x} = \mathbf{x}^\top \left( \frac{1}{2} \mathbf{Q} + \frac{1}{2} \mathbf{Q}^\top \right) \mathbf{x}.$$

A quadratic form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$ ,  $\mathbf{Q} = \mathbf{Q}^\top$ , is said to be *positive definite* if  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} > 0$  for all nonzero vectors  $\mathbf{x}$ . It is *positive semidefinite* if  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} \geq 0$  for all  $\mathbf{x}$ . Similarly, we define the quadratic form to be *negative definite*, or *negative semidefinite*, if  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} < 0$  for all nonzero vectors  $\mathbf{x}$ , or  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} \leq 0$  for all  $\mathbf{x}$ , respectively.

Recall that the minors of a matrix  $\mathbf{Q}$  are the determinants of the matrices obtained by successively removing rows and columns from  $\mathbf{Q}$ . The *principal minors* are  $\det \mathbf{Q}$  itself and the determinants of matrices obtained by successively removing an  $i$ th row and an  $i$ th column. That is, the principal minors are

$$\det \begin{bmatrix} q_{i_1 i_1} & q_{i_1 i_2} & \cdots & q_{i_1 i_p} \\ q_{i_2 i_1} & q_{i_2 i_2} & \cdots & q_{i_2 i_p} \\ \vdots & \vdots & & \vdots \\ q_{i_p i_1} & q_{i_p i_2} & \cdots & q_{i_p i_p} \end{bmatrix}, \quad 1 \leq i_1 < \cdots < i_p \leq n, \quad p = 1, 2, \dots, n.$$

The *leading principal minors* are  $\det \mathbf{Q}$  and the minors obtained by successively removing the last row and the last column. That is, the leading principal minors are

$$\begin{aligned} \Delta_1 &= q_{11}, & \Delta_2 &= \det \begin{bmatrix} q_{11} & q_{12} \\ q_{21} & q_{22} \end{bmatrix}, \\ \Delta_3 &= \det \begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix}, & \dots, & \Delta_n = \det \mathbf{Q}. \end{aligned}$$

We now prove *Sylvester's criterion*, which allows us to determine if a quadratic form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$  is positive definite using only the leading principal minors of  $\mathbf{Q}$ .

**Theorem 3.6 Sylvester's Criterion.** A quadratic form  $\mathbf{x}^\top \mathbf{Q}\mathbf{x}$ ,  $\mathbf{Q} = \mathbf{Q}^\top$ , is positive definite if and only if the leading principal minors of  $\mathbf{Q}$  are positive.

□

*Proof.* The key to the proof of Sylvester's criterion is the fact that a quadratic form whose leading principal minors are nonzero can be expressed in some basis as a sum of squares

$$\frac{\Delta_0}{\Delta_1} \tilde{x}_1^2 + \frac{\Delta_1}{\Delta_2} \tilde{x}_2^2 + \cdots + \frac{\Delta_{n-1}}{\Delta_n} \tilde{x}_n^2,$$

where  $\tilde{x}_i$  are the coordinates of the vector  $\mathbf{x}$  in the new basis,  $\Delta_0 \triangleq 1$ , and  $\Delta_1, \dots, \Delta_n$  are the leading principal minors of  $\mathbf{Q}$ .

To this end, consider a quadratic form  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x}$ , where  $\mathbf{Q} = \mathbf{Q}^\top$ . Let  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$  be the natural basis for  $\mathbb{R}^n$ , and let

$$\mathbf{x} = x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 + \cdots + x_n \mathbf{e}_n$$

be a given vector in  $\mathbb{R}^n$ . Let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  be another basis for  $\mathbb{R}^n$ . Then, the vector  $\mathbf{x}$  is represented in the new basis as  $\tilde{\mathbf{x}}$ , where

$$\mathbf{x} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n] \tilde{\mathbf{x}} \triangleq \mathbf{V} \tilde{\mathbf{x}}.$$

Accordingly, the quadratic form can be written as

$$\mathbf{x}^\top \mathbf{Q}\mathbf{x} = \tilde{\mathbf{x}}^\top \mathbf{V}^\top \mathbf{Q} \mathbf{V} \tilde{\mathbf{x}} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}},$$

where

$$\tilde{\mathbf{Q}} = \mathbf{V}^\top \mathbf{Q} \mathbf{V} = \begin{bmatrix} \tilde{q}_{11} & \cdots & \tilde{q}_{1n} \\ \vdots & \ddots & \vdots \\ \tilde{q}_{n1} & \cdots & \tilde{q}_{nn} \end{bmatrix}.$$

Note that  $\tilde{q}_{ij} = \langle \mathbf{v}_i, \mathbf{Q}\mathbf{v}_j \rangle$ . Our goal is to determine conditions on the new basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  such that  $\tilde{q}_{ij} = 0$  for  $i \neq j$ .

We seek the new basis in the form

$$\begin{aligned} \mathbf{v}_1 &= \alpha_{11} \mathbf{e}_1, \\ \mathbf{v}_2 &= \alpha_{21} \mathbf{e}_1 + \alpha_{22} \mathbf{e}_2, \\ &\vdots \\ \mathbf{v}_n &= \alpha_{n1} \mathbf{e}_1 + \alpha_{n2} \mathbf{e}_2 + \cdots + \alpha_{nn} \mathbf{e}_n. \end{aligned}$$

Observe that for  $j = 1, \dots, i-1$ , if

$$\langle \mathbf{v}_i, \mathbf{Q}\mathbf{e}_j \rangle = 0,$$

then

$$\langle \mathbf{v}_i, \mathbf{Q}\mathbf{v}_j \rangle = 0.$$

Our goal then is to determine the coefficients  $\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{ii}$ ,  $i = 1, \dots, n$ , such that the vector

$$\mathbf{v}_i = \alpha_{i1}\mathbf{e}_1 + \alpha_{i2}\mathbf{e}_2 + \cdots + \alpha_{ii}\mathbf{e}_i$$

satisfies the  $i$  relations

$$\begin{aligned} \langle \mathbf{v}_i, \mathbf{Q}\mathbf{e}_j \rangle &= 0, \quad j = 1, \dots, i-1, \\ \langle \mathbf{e}_i, \mathbf{Q}\mathbf{v}_i \rangle &= 1. \end{aligned}$$

In this case, we get

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \alpha_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \alpha_{nn} \end{bmatrix}.$$

For each  $i = 1, \dots, n$ , the  $i$  relations above determine the coefficients  $\alpha_{i1}, \dots, \alpha_{ii}$  in a unique way. Indeed, upon substituting the expression for  $\mathbf{v}_i$  into the equations above, we obtain the set of equations

$$\alpha_{i1}q_{11} + \alpha_{i2}q_{12} + \cdots + \alpha_{ii}q_{1i} = 0,$$

$$\vdots$$

$$\alpha_{i1}q_{i-11} + \alpha_{i2}q_{i-12} + \cdots + \alpha_{ii}q_{i-1i} = 0,$$

$$\alpha_{i1}q_{ii} + \alpha_{i2}q_{i2} + \cdots + \alpha_{ii}q_{ii} = 1.$$

The set of equations above can be expressed in matrix form as

$$\begin{bmatrix} q_{11} & q_{12} & \cdots & q_{1i} \\ q_{21} & q_{22} & \cdots & q_{2i} \\ \vdots & \vdots & \ddots & \vdots \\ q_{i1} & q_{i2} & \cdots & q_{ii} \end{bmatrix} \begin{bmatrix} \alpha_{i1} \\ \alpha_{i2} \\ \vdots \\ \alpha_{ii} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

If the leading principal minors of the matrix  $\mathbf{Q}$  do not vanish, then the coefficients  $\alpha_{ij}$  can be obtained using *Cramer's rule*. In particular,

$$\alpha_{ii} = \frac{1}{\Delta_i} \det \begin{bmatrix} q_{11} & \cdots & q_{1i-1} & 0 \\ \vdots & \ddots & \vdots & 0 \\ q_{i-11} & \cdots & q_{i-1i-1} & 0 \\ q_{ii} & \cdots & q_{ii-1} & 1 \end{bmatrix} = \frac{\Delta_{i-1}}{\Delta_i}.$$

Hence,

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \frac{1}{\Delta_1} & & & 0 \\ & \frac{\Delta_1}{\Delta_2} & & \\ & & \ddots & \\ 0 & & & \frac{\Delta_{n-1}}{\Delta_n} \end{bmatrix}.$$

In the new basis, the quadratic form can be expressed as a sum of squares

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}} = \frac{1}{\Delta_1} \tilde{x}_1^2 + \frac{\Delta_1}{\Delta_2} \tilde{x}_2^2 + \cdots + \frac{\Delta_{n-1}}{\Delta_n} \tilde{x}_n^2.$$

We now show that a necessary and sufficient condition for the quadratic form to be positive definite is  $\Delta_i > 0$ ,  $i = 1, \dots, n$ .

Sufficiency is clear, for if  $\Delta_i > 0$ ,  $i = 1, \dots, n$ , then by the previous argument there is a basis such that

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}} > 0$$

for any  $\mathbf{x} \neq \mathbf{0}$  (or, equivalently, any  $\tilde{\mathbf{x}} \neq \mathbf{0}$ ).

To prove necessity, we first show that for  $i = 1, \dots, n$ , we have  $\Delta_i \neq 0$ . To see this, suppose that  $\Delta_k = 0$  for some  $k$ . Note that  $\Delta_k = \det \mathbf{Q}_k$ ,

$$\mathbf{Q}_k = \begin{bmatrix} q_{11} & \cdots & q_{1k} \\ \vdots & \ddots & \vdots \\ q_{k1} & \cdots & q_{kk} \end{bmatrix}.$$

Then, there exists a vector  $\mathbf{v} \in \mathbb{R}^k$ ,  $\mathbf{v} \neq \mathbf{0}$ , such that  $\mathbf{v}^\top \mathbf{Q}_k = \mathbf{0}$ . Now let  $\mathbf{x} \in \mathbb{R}^n$  be given by  $\mathbf{x} = [\mathbf{v}^\top, \mathbf{0}^\top]^\top$ . Then,

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \mathbf{v}^\top \mathbf{Q}_k \mathbf{v} = 0.$$

But  $\mathbf{x} \neq \mathbf{0}$ , which contradicts the fact that the quadratic form  $f$  is positive definite. Therefore, if  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} > 0$ , then  $\Delta_i \neq 0$ ,  $i = 1, \dots, n$ . Then, using our previous argument, we may write

$$\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{Q}} \tilde{\mathbf{x}} = \frac{1}{\Delta_1} \tilde{x}_1^2 + \frac{\Delta_1}{\Delta_2} \tilde{x}_2^2 + \cdots + \frac{\Delta_{n-1}}{\Delta_n} \tilde{x}_n^2,$$

where  $\tilde{\mathbf{x}} = [\mathbf{v}_1, \dots, \mathbf{v}_n] \mathbf{x}$ . Hence, if the quadratic form is positive definite, then all leading principal minors must be positive. ■

Note that if  $\mathbf{Q}$  is not symmetric, Sylvester's criterion cannot be used to check positive definiteness of the quadratic form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$ . To see this, consider an example where

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 \\ -4 & 1 \end{bmatrix}.$$

The leading principal minors of  $\mathbf{Q}$  are  $\Delta_1 = 1 > 0$  and  $\Delta_2 = \det \mathbf{Q} = 1 > 0$ . However, if  $\mathbf{x} = [1, 1]^\top$ , then  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} = -2 < 0$ , and hence the associated quadratic form is not positive definite. Note that

$$\begin{aligned} \mathbf{x}^\top \mathbf{Q} \mathbf{x} &= \mathbf{x}^\top \begin{bmatrix} 1 & 0 \\ -4 & 1 \end{bmatrix} \mathbf{x} = \frac{1}{2} \mathbf{x}^\top \left( \begin{bmatrix} 1 & 0 \\ -4 & 1 \end{bmatrix} + \begin{bmatrix} 1 & -4 \\ 0 & 1 \end{bmatrix} \right) \mathbf{x} \\ &= \mathbf{x}^\top \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \mathbf{x} = \mathbf{x}^\top \mathbf{Q}_0 \mathbf{x}. \end{aligned}$$

The leading principal minors of  $\mathbf{Q}_0$  are  $\Delta_1 = 1 > 0$  and  $\Delta_2 = \det \mathbf{Q}_0 = -3 < 0$ , as expected.

A necessary condition for a real quadratic form to be positive semidefinite is that the leading principal minors be nonnegative. However, this is *not* a sufficient condition (see Exercise 3.16). In fact, a real quadratic form is positive semidefinite if and only if all principal minors are nonnegative (for a proof of this fact, see [44, p. 307]).

A symmetric matrix  $\mathbf{Q}$  is said to be *positive definite* if the quadratic form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$  is positive definite. If  $\mathbf{Q}$  is positive definite, we write  $\mathbf{Q} > 0$ . Similarly, we define a symmetric matrix  $\mathbf{Q}$  to be *positive semidefinite* ( $\mathbf{Q} \geq 0$ ), *negative definite* ( $\mathbf{Q} < 0$ ), and *negative semidefinite* ( $\mathbf{Q} \leq 0$ ) if the corresponding quadratic forms have the respective properties. The symmetric matrix  $\mathbf{Q}$  is *indefinite* if it is neither positive semidefinite nor negative semidefinite. Note that the matrix  $\mathbf{Q}$  is positive definite (semidefinite) if and only if the matrix  $-\mathbf{Q}$  is negative definite (semidefinite).

Sylvester's criterion provides a way of checking the definiteness of a quadratic form, or equivalently, a symmetric matrix. An alternative method involves checking the eigenvalues of  $\mathbf{Q}$ , as stated below.

**Theorem 3.7** *A symmetric matrix  $\mathbf{Q}$  is positive definite (or positive semidefinite) if and only if all eigenvalues of  $\mathbf{Q}$  are positive (or nonnegative).*  $\square$

*Proof.* For any  $\mathbf{x}$ , let  $\mathbf{y} = \mathbf{T}^{-1}\mathbf{x} = \mathbf{T}^\top \mathbf{x}$ , where  $\mathbf{T}$  is an orthogonal matrix whose columns are eigenvectors of  $\mathbf{Q}$ . Then,  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} = \mathbf{y}^\top \mathbf{T}^\top \mathbf{Q} \mathbf{T} \mathbf{y} = \sum_{i=1}^n \lambda_i y_i^2$ . From this, the result follows.  $\blacksquare$

Through diagonalization, we can show that a symmetric positive semidefinite matrix  $\mathbf{Q}$  has a positive semidefinite (symmetric) square root  $\mathbf{Q}^{1/2}$  satisfying  $\mathbf{Q}^{1/2} \mathbf{Q}^{1/2} = \mathbf{Q}$ . For this, we use  $\mathbf{T}$  as above and define

$$\mathbf{Q}^{1/2} = \mathbf{T} \begin{bmatrix} \lambda_1^{1/2} & & & 0 \\ & \lambda_2^{1/2} & & \\ & & \ddots & \\ 0 & & & \lambda_n^{1/2} \end{bmatrix} \mathbf{T}^\top,$$

which is easily verified to have the desired properties. Note that the quadratic form  $\mathbf{x}^\top \mathbf{Q} \mathbf{x}$  can be expressed as  $\|\mathbf{Q}^{1/2} \mathbf{x}\|^2$ .

In summary, we have presented two tests for definiteness of quadratic forms and symmetric matrices. We point out again that nonnegativity of leading principal minors is a necessary but not a sufficient condition for positive semidefiniteness.

### 3.5 Matrix Norms

The norm of a matrix may be chosen in a variety of ways. Because the set of matrices  $\mathbb{R}^{m \times n}$  can be viewed as the real vector space  $\mathbb{R}^{mn}$ , matrix norms

should be no different from regular vector norms. Therefore, we define the norm of a matrix  $\mathbf{A}$ , denoted  $\|\mathbf{A}\|$ , to be any function  $\|\cdot\|$  that satisfies the following conditions:

1.  $\|\mathbf{A}\| > 0$  if  $\mathbf{A} \neq \mathbf{O}$ , and  $\|\mathbf{O}\| = 0$ , where  $\mathbf{O}$  is a matrix with all entries equal to zero.
2.  $\|c\mathbf{A}\| = |c|\|\mathbf{A}\|$ , for any  $c \in \mathbb{R}$ .
3.  $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ .

An example of a matrix norm is the *Frobenius norm*, defined as

$$\|\mathbf{A}\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n (a_{ij})^2 \right)^{\frac{1}{2}},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Note that the Frobenius norm is equivalent to the Euclidean norm on  $\mathbb{R}^{mn}$ .

For our purposes, we consider only matrix norms that satisfy the following additional condition:

4.  $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$ .

It turns out that the Frobenius norm satisfies condition 4 as well.

In many problems, both matrices and vectors appear simultaneously. Therefore, it is convenient to construct the norm of a matrix in such a way that it will be related to vector norms. To this end we consider a special class of matrix norms, called *induced norms*. Let  $\|\cdot\|_{(n)}$  and  $\|\cdot\|_{(m)}$  be vector norms on  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , respectively. We say that the matrix norm is *induced* by, or is *compatible* with, the given vector norms if for any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and any vector  $\mathbf{x} \in \mathbb{R}^n$ , the following inequality is satisfied:

$$\|\mathbf{Ax}\|_{(m)} \leq \|\mathbf{A}\| \|\mathbf{x}\|_{(n)}.$$

We can define an induced matrix norm as

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|_{(n)}=1} \|\mathbf{Ax}\|_{(m)};$$

that is,  $\|\mathbf{A}\|$  is the maximum of the norms of the vectors  $\mathbf{Ax}$  where the vector  $\mathbf{x}$  runs over the set of all vectors with unit norm. When there is no ambiguity, we omit the subscripts  $(m)$  and  $(n)$  from  $\|\cdot\|_{(m)}$  and  $\|\cdot\|_{(n)}$ .

Because of the continuity of a vector norm (see Exercise 2.10), for each matrix  $\mathbf{A}$  the maximum

$$\max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|$$

is attainable; that is, a vector  $\mathbf{x}_0$  exists such that  $\|\mathbf{x}_0\| = 1$  and  $\|\mathbf{Ax}_0\| = \|\mathbf{A}\|$ . This fact follows from the theorem of Weierstrass (see Theorem 4.2).

The induced norm satisfies conditions 1 to 4 and the compatibility condition, as we prove below.

*Proof of Condition 1.* Let  $\mathbf{A} \neq \mathbf{O}$ . Then, a vector  $\mathbf{x}$ ,  $\|\mathbf{x}\| = 1$ , can be found such that  $\mathbf{Ax} \neq \mathbf{0}$ , and thus  $\|\mathbf{Ax}\| \neq 0$ . Hence,  $\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\| \neq 0$ . If, on the other hand,  $\mathbf{A} = \mathbf{O}$ , then  $\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ox}\| = 0$ . ■

*Proof of Condition 2.* By definition,  $\|c\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|c\mathbf{Ax}\|$ . Obviously,  $\|c\mathbf{Ax}\| = |c|\|\mathbf{Ax}\|$ , and therefore  $\|c\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} |c|\|\mathbf{Ax}\| = |c| \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\| = |c|\|\mathbf{A}\|$ . ■

*Proof of Compatibility Condition.* Let  $\mathbf{y} \neq \mathbf{0}$  be any vector. Then,  $\mathbf{x} = \mathbf{y}/\|\mathbf{y}\|$  satisfies the condition  $\|\mathbf{x}\| = 1$ . Consequently,  $\|\mathbf{Ay}\| = \|\mathbf{A}(\|\mathbf{y}\|\mathbf{x})\| = \|\mathbf{y}\|\|\mathbf{Ax}\| \leq \|\mathbf{y}\|\|\mathbf{A}\|$ . ■

*Proof of Condition 3.* For the matrix  $\mathbf{A} + \mathbf{B}$ , we can find a vector  $\mathbf{x}_0$  such that  $\|\mathbf{A} + \mathbf{B}\| = \|(\mathbf{A} + \mathbf{B})\mathbf{x}_0\|$  and  $\|\mathbf{x}_0\| = 1$ . Then, we have

$$\begin{aligned}\|\mathbf{A} + \mathbf{B}\| &= \|(\mathbf{A} + \mathbf{B})\mathbf{x}_0\| \\ &= \|\mathbf{Ax}_0 + \mathbf{Bx}_0\| \\ &\leq \|\mathbf{Ax}_0\| + \|\mathbf{Bx}_0\| \\ &\leq \|\mathbf{A}\|\|\mathbf{x}_0\| + \|\mathbf{B}\|\|\mathbf{x}_0\| \\ &= \|\mathbf{A}\| + \|\mathbf{B}\|,\end{aligned}$$

which shows that condition 3 holds. ■

*Proof of Condition 4.* For the matrix  $\mathbf{AB}$ , we can find a vector  $\mathbf{x}_0$  such that  $\|\mathbf{x}_0\| = 1$  and  $\|\mathbf{ABx}_0\| = \|\mathbf{AB}\|$ . Then, we have

$$\begin{aligned}\|\mathbf{AB}\| &= \|\mathbf{ABx}_0\| \\ &= \|\mathbf{A}(\mathbf{Bx}_0)\| \\ &\leq \|\mathbf{A}\|\|\mathbf{Bx}_0\| \\ &\leq \|\mathbf{A}\|\|\mathbf{B}\|\|\mathbf{x}_0\| \\ &= \|\mathbf{A}\|\|\mathbf{B}\|,\end{aligned}$$

which shows that condition 4 holds. ■

**Theorem 3.8** *Let*

$$\|\mathbf{x}\| = \left( \sum_{k=1}^n |x_k|^2 \right)^{1/2} = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

*The matrix norm induced by this vector norm is*

$$\|\mathbf{A}\| = \sqrt{\lambda_1},$$

where  $\lambda_1$  is the largest eigenvalue of the matrix  $\mathbf{A}^\top \mathbf{A}$ .  $\square$

*Proof.* We have

$$\|\mathbf{A}\mathbf{x}\|^2 = \langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{A}^\top \mathbf{A}\mathbf{x} \rangle.$$

The matrix  $\mathbf{A}^\top \mathbf{A}$  is symmetric and positive semidefinite. Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$  be its eigenvalues and  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  the orthonormal set of the eigenvectors corresponding to these eigenvalues. Now, we take an arbitrary vector  $\mathbf{x}$  with  $\|\mathbf{x}\| = 1$  and represent it as a linear combination of  $\mathbf{x}_i$ ,  $i = 1, \dots, n$ :

$$\mathbf{x} = c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_n \mathbf{x}_n.$$

Note that

$$\langle \mathbf{x}, \mathbf{x} \rangle = c_1^2 + c_2^2 + \dots + c_n^2 = 1.$$

Furthermore,

$$\begin{aligned} \|\mathbf{A}\mathbf{x}\|^2 &= \langle \mathbf{x}, \mathbf{A}^\top \mathbf{A}\mathbf{x} \rangle \\ &= \langle c_1 \mathbf{x}_1 + \dots + c_n \mathbf{x}_n, c_1 \lambda_1 \mathbf{x}_1 + \dots + c_n \lambda_n \mathbf{x}_n \rangle \\ &= \lambda_1 c_1^2 + \dots + \lambda_n c_n^2 \\ &\leq \lambda_1 (c_1^2 + \dots + c_n^2) \\ &= \lambda_1. \end{aligned}$$

For the eigenvector  $\mathbf{x}_1$  of  $\mathbf{A}^\top \mathbf{A}$  corresponding to the eigenvalue  $\lambda_1$ , we have

$$\|\mathbf{A}\mathbf{x}_1\|^2 = \langle \mathbf{x}_1, \mathbf{A}^\top \mathbf{A}\mathbf{x}_1 \rangle = \langle \mathbf{x}_1, \lambda_1 \mathbf{x}_1 \rangle = \lambda_1,$$

and hence

$$\max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| = \sqrt{\lambda_1}.$$

This completes the proof.  $\blacksquare$

Using arguments similar to the above, we can deduce the following important inequalities.

**Rayleigh's Inequalities.** If an  $n \times n$  matrix  $\mathbf{P}$  is real symmetric positive definite, then

$$\lambda_{\min}(\mathbf{P})\|\mathbf{x}\|^2 \leq \mathbf{x}^\top \mathbf{P}\mathbf{x} \leq \lambda_{\max}(\mathbf{P})\|\mathbf{x}\|^2,$$

where  $\lambda_{\min}(\mathbf{P})$  denotes the smallest eigenvalue of  $\mathbf{P}$ , and  $\lambda_{\max}(\mathbf{P})$  denotes the largest eigenvalue of  $\mathbf{P}$ .

**Example 3.1** Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix},$$

and let the norm in  $\mathbb{R}^2$  be given by

$$\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2}.$$

Then,

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix},$$

and  $\det[\lambda \mathbf{I}_2 - \mathbf{A}^\top \mathbf{A}] = \lambda^2 - 10\lambda + 9 = (\lambda - 1)(\lambda - 9)$ . Thus,  $\|\mathbf{A}\| = \sqrt{9} = 3$ .

The eigenvector of  $\mathbf{A}^\top \mathbf{A}$  corresponding to  $\lambda_1 = 9$  is

$$\mathbf{x}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

Note that  $\|\mathbf{A}\mathbf{x}_1\| = \|\mathbf{A}\|$ . Indeed,

$$\begin{aligned} \|\mathbf{A}\mathbf{x}_1\| &= \left\| \frac{1}{\sqrt{2}} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\| \\ &= \frac{1}{\sqrt{2}} \left\| \begin{bmatrix} 3 \\ 3 \end{bmatrix} \right\| \\ &= \frac{1}{\sqrt{2}} \sqrt{3^2 + 3^2} \\ &= 3. \end{aligned}$$

Because  $\mathbf{A} = \mathbf{A}^\top$  in this example, we also have  $\|\mathbf{A}\| = \max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})|$ , where  $\lambda_1(\mathbf{A}), \dots, \lambda_n(\mathbf{A})$  are the eigenvalues of  $\mathbf{A}$  (possibly repeated). ■

**Warning:** In general,  $\max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})| \neq \|\mathbf{A}\|$ . Instead, we have  $\|\mathbf{A}\| \geq \max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})|$ , as illustrated in the following example (see also Exercise 5.2).

**Example 3.2** Let

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

then

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

and

$$\det[\lambda \mathbf{I}_2 - \mathbf{A}^\top \mathbf{A}] = \det \begin{bmatrix} \lambda & 0 \\ 0 & \lambda - 1 \end{bmatrix} = \lambda(\lambda - 1).$$

Note that 0 is the only eigenvalue of  $\mathbf{A}$ . Thus, for  $i = 1, 2$ ,  $\|\mathbf{A}\| = 1 > |\lambda_i(\mathbf{A})| = 0$ . ■

For a more complete but still basic treatment of topics in linear algebra as discussed in this and the preceding chapter, see [47], [66], [95], [126]. For a treatment of matrices, we refer the reader to [44], [62]. Numerical aspects of matrix computations are discussed in [41], [53].

## EXERCISES

**3.1** Fix a vector in  $\mathbb{R}^n$ . Let  $\mathbf{x}$  be the column of the coordinates of the vector with respect to the basis  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$  and  $\mathbf{x}'$  the coordinates of the same vector with respect to the basis  $\{\mathbf{e}'_1, \mathbf{e}'_2, \dots, \mathbf{e}'_n\}$ . Show that  $\mathbf{x}' = \mathbf{T}\mathbf{x}$ , where  $\mathbf{T}$  is the transformation matrix from  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$  to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \dots, \mathbf{e}'_n\}$ .

**3.2** For each of the following cases, find the transformation matrix  $\mathbf{T}$  from  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  to  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$ :

- a.  $\mathbf{e}'_1 = \mathbf{e}_1 + 3\mathbf{e}_2 - 4\mathbf{e}_3$ ,  $\mathbf{e}'_2 = 2\mathbf{e}_1 - \mathbf{e}_2 + 5\mathbf{e}_3$ ,  $\mathbf{e}'_3 = 4\mathbf{e}_1 + 5\mathbf{e}_2 + 3\mathbf{e}_3$ .
- b.  $\mathbf{e}_1 = \mathbf{e}'_1 + \mathbf{e}'_2 + 3\mathbf{e}'_3$ ,  $\mathbf{e}_2 = 2\mathbf{e}'_1 - \mathbf{e}'_2 + 4\mathbf{e}'_3$ ,  $\mathbf{e}_3 = 3\mathbf{e}'_1 + 5\mathbf{e}'_3$ .

**3.3** Consider two bases of  $\mathbb{R}^3$ ,  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  and  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$ , where  $\mathbf{e}_1 = 2\mathbf{e}'_1 + \mathbf{e}'_2 - \mathbf{e}'_3$ ,  $\mathbf{e}_2 = 2\mathbf{e}'_1 - \mathbf{e}'_2 + 2\mathbf{e}'_3$ , and  $\mathbf{e}_3 = 3\mathbf{e}'_1 + \mathbf{e}'_3$ . Suppose that a linear transformation has a matrix representation in  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  of the form

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{bmatrix}.$$

Find the matrix representation of this linear transformation in the basis  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$ .

**3.4** Consider two bases of  $\mathbb{R}^4$ ,  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4\}$  and  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3, \mathbf{e}'_4\}$ , where  $\mathbf{e}'_1 = \mathbf{e}_1$ ,  $\mathbf{e}'_2 = \mathbf{e}_1 + \mathbf{e}_2$ ,  $\mathbf{e}'_3 = \mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3$ , and  $\mathbf{e}'_4 = \mathbf{e}_1 + \mathbf{e}_2 + \mathbf{e}_3 + \mathbf{e}_4$ . Suppose that a linear transformation has a matrix representation in  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4\}$  of the form

$$\mathbf{A} = \begin{bmatrix} 2 & 0 & 1 & 0 \\ -3 & 2 & 0 & 1 \\ 0 & 1 & -1 & 2 \\ 1 & 0 & 0 & 3 \end{bmatrix}.$$

Find the matrix representation of this linear transformation in the basis  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3, \mathbf{e}'_4\}$ .

**3.5** Consider a linear transformation given by the matrix

$$\mathbf{A} = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 2 & 5 & 2 & 1 \\ -1 & 1 & 0 & 3 \end{bmatrix}.$$

Find a basis for  $\mathbb{R}^4$  with respect to which the matrix representation for the linear transformation above is diagonal.

**3.6** Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of the matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Show that the eigenvalues of the matrix  $\mathbf{I}_n - \mathbf{A}$  are  $1 - \lambda_1, \dots, 1 - \lambda_n$ .

**3.7** Let  $\mathcal{V}$  be a subspace. Show that  $\mathcal{V}^\perp$  is also a subspace.

**3.8** Find the nullspace of

$$\mathbf{A} = \begin{bmatrix} 4 & -2 & 0 \\ 2 & 1 & -1 \\ 2 & -3 & 1 \end{bmatrix}.$$

**3.9** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a matrix. Show that  $\mathcal{R}(\mathbf{A})$  is a subspace of  $\mathbb{R}^m$  and  $\mathcal{N}(\mathbf{A})$  is a subspace of  $\mathbb{R}^n$ .

**3.10** Prove that if  $\mathbf{A}$  and  $\mathbf{B}$  are two matrices with  $m$  rows, and  $\mathcal{N}(\mathbf{A}^\top) \subset \mathcal{N}(\mathbf{B}^\top)$ , then  $\mathcal{R}(\mathbf{B}) \subset \mathcal{R}(\mathbf{A})$ .

*Hint:* Use the fact that for any matrix  $\mathbf{M}$  with  $m$  rows, we have  $\dim \mathcal{R}(\mathbf{M}) + \dim \mathcal{N}(\mathbf{M}^\top) = m$  [this is one of the fundamental theorems of linear algebra (see [126, p. 75])].

**3.11** Let  $\mathcal{V}$  be a subspace. Show that  $(\mathcal{V}^\perp)^\perp = \mathcal{V}$ .

*Hint:* Use Exercise 3.10.

**3.12** Let  $\mathcal{V}$  and  $\mathcal{W}$  be subspaces. Show that if  $\mathcal{V} \subset \mathcal{W}$ , then  $\mathcal{W}^\perp \subset \mathcal{V}^\perp$ .

**3.13** Let  $\mathcal{V}$  be a subspace of  $\mathbb{R}^n$ . Show that there exist matrices  $\mathbf{V}$  and  $\mathbf{U}$  such that  $\mathcal{V} = \mathcal{R}(\mathbf{V}) = \mathcal{N}(\mathbf{U})$ .

**3.14** Let  $\mathbf{P}$  be an orthogonal projector onto a subspace  $\mathcal{V}$ . Show that

a.  $\mathbf{P}\mathbf{x} = \mathbf{x}$  for all  $\mathbf{x} \in \mathcal{V}$ .

b.  $\mathcal{R}(\mathbf{P}) = \mathcal{V}$ .

**3.15** Is the quadratic form

$$\mathbf{x}^\top \begin{bmatrix} 1 & -8 \\ 1 & 1 \end{bmatrix} \mathbf{x}$$

positive definite, positive semidefinite, negative definite, negative semidefinite, or indefinite?

**3.16** Let

$$\mathbf{A} = \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 0 \end{bmatrix}.$$

Show that although all leading principal minors of  $\mathbf{A}$  are nonnegative,  $\mathbf{A}$  is not positive semidefinite.

**3.17** Consider the matrix

$$\mathbf{Q} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}.$$

- a. Is this matrix positive definite, negative definite, or indefinite?
- b. Is this matrix positive definite, negative definite, or indefinite on the subspace

$$\mathcal{M} = \{\mathbf{x} : x_1 + x_2 + x_3 = 0\} ?$$

**3.18** For each of the following quadratic forms, determine if it is positive definite, negative definite, positive semidefinite, negative semidefinite, or indefinite.

- a.  $f(x_1, x_2, x_3) = x_2^2$
- b.  $f(x_1, x_2, x_3) = x_1^2 + 2x_2^2 - x_1x_3$
- c.  $f(x_1, x_2, x_3) = x_1^2 + x_3^2 + 2x_1x_2 + 2x_1x_3 + 2x_2x_3$

**3.19** Find a transformation that brings the following quadratic form into the diagonal form:

$$f(x_1, x_2, x_3) = 4x_1^2 + x_2^2 + 9x_3^2 - 4x_1x_2 - 6x_2x_3 + 12x_1x_3.$$

*Hint:* Read the proof of Theorem 3.6.

**3.20** Consider the quadratic form

$$f(x_1, x_2, x_3) = x_1^2 + x_2^2 + 5x_3^2 + 2\xi x_1 x_2 - 2x_1 x_3 + 4x_2 x_3.$$

Find the values of the parameter  $\xi$  for which this quadratic form is positive definite.

**3.21** Consider the function  $\langle \cdot, \cdot \rangle_Q : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , defined by  $\langle \mathbf{x}, \mathbf{y} \rangle_Q = \mathbf{x}^\top \mathbf{Q} \mathbf{y}$ , where  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a symmetric positive definite matrix. Show that  $\langle \cdot, \cdot \rangle_Q$  satisfies conditions 1 to 4 for inner products (see Section 2.4).

**3.22** Consider the vector norm  $\|\cdot\|_\infty$  on  $\mathbb{R}^n$  given by  $\|\mathbf{x}\|_\infty = \max_i |x_i|$ , where  $\mathbf{x} = [x_1, \dots, x_n]^\top$ . Define the norm  $\|\cdot\|_\infty$  on  $\mathbb{R}^m$  similarly. Show that the matrix norm induced by these vector norms is given by

$$\|\mathbf{A}\|_\infty = \max_i \sum_{k=1}^n |a_{ik}|,$$

where  $a_{ij}$  is the  $(i, j)$ th element of  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .

**3.23** Consider the vector norm  $\|\cdot\|_1$  on  $\mathbb{R}^n$  given by  $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ , where  $\mathbf{x} = [x_1, \dots, x_n]^\top$ . Define the norm  $\|\cdot\|_1$  on  $\mathbb{R}^m$  similarly. Show that the matrix norm induced by these vector norms is given by

$$\|\mathbf{A}\|_1 = \max_k \sum_{i=1}^m |a_{ik}|,$$

where  $a_{ij}$  is the  $(i, j)$ th element of  $\mathbf{A} \in \mathbb{R}^{m \times n}$ .



# CHAPTER 4

---

## CONCEPTS FROM GEOMETRY

---

### 4.1 Line Segments

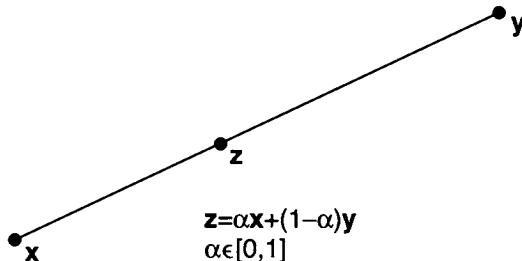
In the following analysis we concern ourselves only with  $\mathbb{R}^n$ . The elements of this space are the  $n$ -component vectors  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$ .

The *line segment* between two points  $\mathbf{x}$  and  $\mathbf{y}$  in  $\mathbb{R}^n$  is the set of points on the straight line joining points  $\mathbf{x}$  and  $\mathbf{y}$  (see Figure 4.1). Note that if  $\mathbf{z}$  lies on the line segment between  $\mathbf{x}$  and  $\mathbf{y}$ , then

$$\mathbf{z} - \mathbf{y} = \alpha(\mathbf{x} - \mathbf{y}),$$

where  $\alpha$  is a real number from the interval  $[0, 1]$ . The equation above can be rewritten as  $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$ . Hence, the line segment between  $\mathbf{x}$  and  $\mathbf{y}$  can be represented as

$$\{\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} : \alpha \in [0, 1]\}.$$



**Figure 4.1** Line segment.

## 4.2 Hyperplanes and Linear Varieties

Let  $u_1, u_2, \dots, u_n, v \in \mathbb{R}$ , where at least one of the  $u_i$  is nonzero. The set of all points  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  that satisfy the linear equation

$$u_1x_1 + u_2x_2 + \cdots + u_nx_n = v$$

is called a *hyperplane* of the space  $\mathbb{R}^n$ . We may describe the hyperplane by

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} = v\},$$

where

$$\mathbf{u} = [u_1, u_2, \dots, u_n]^\top.$$

A hyperplane is not necessarily a subspace of  $\mathbb{R}^n$  since, in general, it does not contain the origin. For  $n = 2$ , the equation of the hyperplane has the form  $u_1x_1 + u_2x_2 = v$ , which is the equation of a straight line. Thus, straight lines are hyperplanes in  $\mathbb{R}^2$ . In  $\mathbb{R}^3$  (three-dimensional space), hyperplanes are ordinary planes. By translating a hyperplane so that it contains the origin of  $\mathbb{R}^n$ , it becomes a subspace of  $\mathbb{R}^n$  (see Figure 4.2). Because the dimension of this subspace is  $n - 1$ , we say that the hyperplane has dimension  $n - 1$ .

The hyperplane  $H = \{\mathbf{x} : u_1x_1 + \cdots + u_nx_n = v\}$  divides  $\mathbb{R}^n$  into two *half-spaces*. One of these half-spaces consists of the points satisfying the inequality  $u_1x_1 + u_2x_2 + \cdots + u_nx_n \geq v$ , denoted

$$H_+ = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} \geq v\},$$

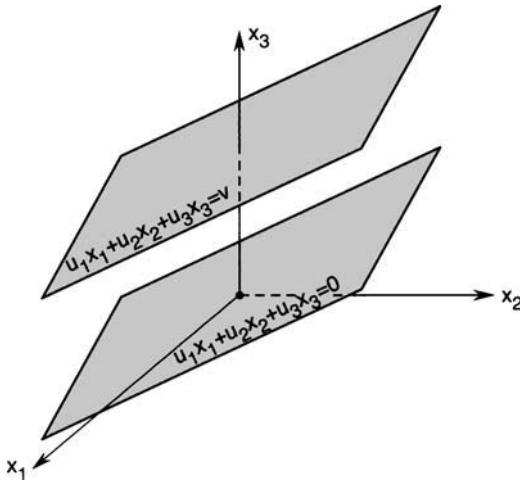
where, as before,

$$\mathbf{u} = [u_1, u_2, \dots, u_n]^\top.$$

The other half-space consists of the points satisfying the inequality  $u_1x_1 + u_2x_2 + \cdots + u_nx_n \leq v$ , denoted

$$H_- = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{u}^\top \mathbf{x} \leq v\}.$$

The half-space  $H_+$  is called the *positive half-space*, and the half-space  $H_-$  is called the *negative half-space*.



**Figure 4.2** Translation of a hyperplane.

Let  $\mathbf{a} = [a_1, a_2, \dots, a_n]^\top$  be an arbitrary point of the hyperplane  $H$ . Thus,  $\mathbf{u}^\top \mathbf{a} - v = 0$ . We can write

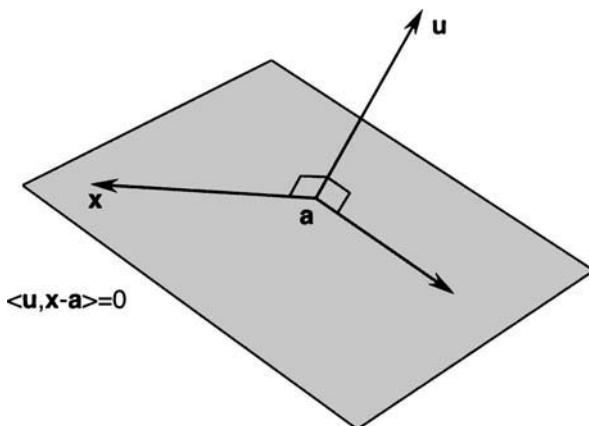
$$\begin{aligned}\mathbf{u}^\top \mathbf{x} - v &= \mathbf{u}^\top \mathbf{x} - v - (\mathbf{u}^\top \mathbf{a} - v) \\ &= \mathbf{u}^\top (\mathbf{x} - \mathbf{a}) \\ &= u_1(x_1 - a_1) + u_2(x_2 - a_2) + \cdots + u_n(x_n - a_n) = 0.\end{aligned}$$

The numbers  $(x_i - a_i)$ ,  $i = 1, \dots, n$ , are the components of the vector  $\mathbf{x} - \mathbf{a}$ . Therefore, the hyperplane  $H$  consists of the points  $\mathbf{x}$  for which  $\langle \mathbf{u}, \mathbf{x} - \mathbf{a} \rangle = 0$ . In other words, the hyperplane  $H$  consists of the points  $\mathbf{x}$  for which the vectors  $\mathbf{u}$  and  $\mathbf{x} - \mathbf{a}$  are orthogonal (see Figure 4.3). We call the vector  $\mathbf{u}$  the *normal* to the hyperplane  $H$ . The set  $H_+$  consists of those points  $\mathbf{x}$  for which  $\langle \mathbf{u}, \mathbf{x} - \mathbf{a} \rangle \geq 0$ , and  $H_-$  consists of those points  $\mathbf{x}$  for which  $\langle \mathbf{u}, \mathbf{x} - \mathbf{a} \rangle \leq 0$ .

A *linear variety* is a set of the form

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$$

for some matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and vector  $\mathbf{b} \in \mathbb{R}^m$ . If  $\dim \mathcal{N}(\mathbf{A}) = r$ , we say that the linear variety has dimension  $r$ . A linear variety is a subspace if and only if  $\mathbf{b} = \mathbf{0}$ . If  $\mathbf{A} = \mathbf{O}$ , the linear variety is  $\mathbb{R}^n$ . If the dimension of the linear variety is less than  $n$ , then it is the intersection of a finite number of hyperplanes.



**Figure 4.3** The hyperplane  $H = \{x \in \mathbb{R}^n : u^\top (x - a) = 0\}$ .

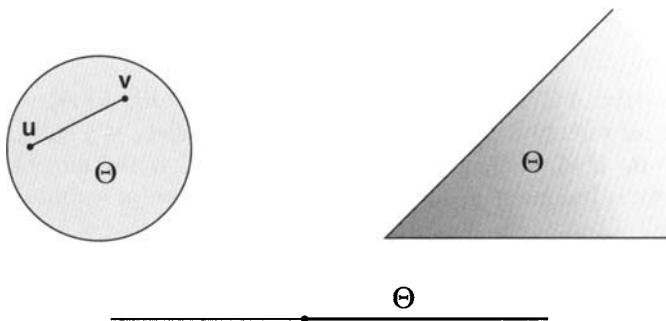
### 4.3 Convex Sets

Recall that the line segment between two points  $u, v \in \mathbb{R}^n$  is the set  $\{w \in \mathbb{R}^n : w = \alpha u + (1 - \alpha)v, \alpha \in [0, 1]\}$ . A point  $w = \alpha u + (1 - \alpha)v$  (where  $\alpha \in [0, 1]$ ) is called a *convex combination* of the points  $u$  and  $v$ .

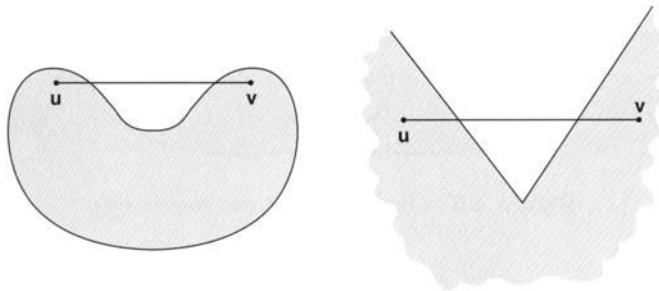
A set  $\Theta \subset \mathbb{R}^n$  is *convex* if for all  $u, v \in \Theta$ , the line segment between  $u$  and  $v$  is in  $\Theta$ . Figure 4.4 gives examples of convex sets, whereas Figure 4.5 gives examples of sets that are not convex. Note that  $\Theta$  is convex if and only if  $\alpha u + (1 - \alpha)v \in \Theta$  for all  $u, v \in \Theta$  and  $\alpha \in (0, 1)$ .

Examples of convex sets include the following:

- The empty set
- A set consisting of a single point
- A line or a line segment



**Figure 4.4** Convex sets.



**Figure 4.5** Sets that are not convex.

- A subspace
- A hyperplane
- A linear variety
- A half-space
- $\mathbb{R}^n$

**Theorem 4.1** Convex subsets of  $\mathbb{R}^n$  have the following properties:

- a. If  $\Theta$  is a convex set and  $\beta$  is a real number, then the set

$$\beta\Theta = \{\mathbf{x} : \mathbf{x} = \beta\mathbf{v}, \mathbf{v} \in \Theta\}$$

is also convex.

- b. If  $\Theta_1$  and  $\Theta_2$  are convex sets, then the set

$$\Theta_1 + \Theta_2 = \{\mathbf{x} : \mathbf{x} = \mathbf{v}_1 + \mathbf{v}_2, \mathbf{v}_1 \in \Theta_1, \mathbf{v}_2 \in \Theta_2\}$$

is also convex.

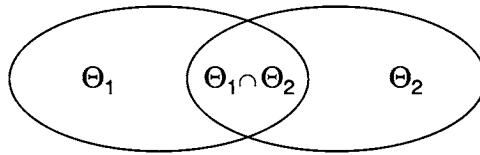
- c. The intersection of any collection of convex sets is convex (see Figure 4.6 for an illustration of this result for two sets).  $\square$

*Proof.*

- a. Let  $\beta\mathbf{v}_1, \beta\mathbf{v}_2 \in \beta\Theta$ , where  $\mathbf{v}_1, \mathbf{v}_2 \in \Theta$ . Because  $\Theta$  is convex, we have  $\alpha\mathbf{v}_1 + (1 - \alpha)\mathbf{v}_2 \in \Theta$  for any  $\alpha \in (0, 1)$ . Hence,

$$\alpha\beta\mathbf{v}_1 + (1 - \alpha)\beta\mathbf{v}_2 = \beta(\alpha\mathbf{v}_1 + (1 - \alpha)\mathbf{v}_2) \in \beta\Theta,$$

and thus  $\beta\Theta$  is convex.



**Figure 4.6** Intersection of two convex sets.

- b. Let  $\mathbf{v}_1, \mathbf{v}_2 \in \Theta_1 + \Theta_2$ . Then,  $\mathbf{v}_1 = \mathbf{v}'_1 + \mathbf{v}''_1$ , and  $\mathbf{v}_2 = \mathbf{v}'_2 + \mathbf{v}''_2$ , where  $\mathbf{v}'_1, \mathbf{v}'_2 \in \Theta_1$ , and  $\mathbf{v}''_1, \mathbf{v}''_2 \in \Theta_2$ . Because  $\Theta_1$  and  $\Theta_2$  are convex, for all  $\alpha \in (0, 1)$ ,

$$\mathbf{x}_1 = \alpha \mathbf{v}'_1 + (1 - \alpha) \mathbf{v}'_2 \in \Theta_1$$

and

$$\mathbf{x}_2 = \alpha \mathbf{v}''_1 + (1 - \alpha) \mathbf{v}''_2 \in \Theta_2.$$

By definition of  $\Theta_1 + \Theta_2$ ,  $\mathbf{x}_1 + \mathbf{x}_2 \in \Theta_1 + \Theta_2$ . Now,

$$\begin{aligned}\alpha \mathbf{v}_1 + (1 - \alpha) \mathbf{v}_2 &= \alpha(\mathbf{v}'_1 + \mathbf{v}''_1) + (1 - \alpha)(\mathbf{v}'_2 + \mathbf{v}''_2) \\ &= \mathbf{x}_1 + \mathbf{x}_2 \in \Theta_1 + \Theta_2.\end{aligned}$$

Hence,  $\Theta_1 + \Theta_2$  is convex.

- c. Let  $C$  be a collection of convex sets. Let  $\mathbf{x}_1, \mathbf{x}_2 \in \bigcap_{\Theta \in C} \Theta$  (where  $\bigcap_{\Theta \in C} \Theta$  represents the intersection of all elements in  $C$ ). Then,  $\mathbf{x}_1, \mathbf{x}_2 \in \Theta$  for each  $\Theta \in C$ . Because each  $\Theta \in C$  is convex,  $\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2 \in \Theta$  for all  $\alpha \in (0, 1)$  and each  $\Theta \in C$ . Thus,  $\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2 \in \bigcap_{\Theta \in C} \Theta$ . ■

A point  $\mathbf{x}$  in a convex set  $\Theta$  is said to be an *extreme point* of  $\Theta$  if there are no two distinct points  $\mathbf{u}$  and  $\mathbf{v}$  in  $\Theta$  such that  $\mathbf{x} = \alpha \mathbf{u} + (1 - \alpha) \mathbf{v}$  for some  $\alpha \in (0, 1)$ . For example, in Figure 4.4, any point on the boundary of the disk is an extreme point, the vertex (corner) of the set on the right is an extreme point, and the endpoint of the half-line is also an extreme point.

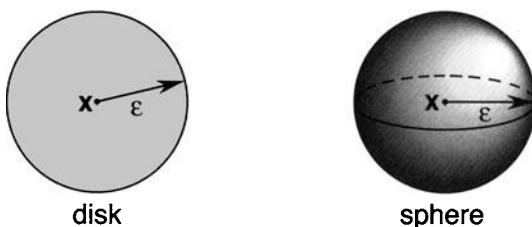
#### 4.4 Neighborhoods

A *neighborhood* of a point  $\mathbf{x} \in \mathbb{R}^n$  is the set

$$\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| < \varepsilon\},$$

where  $\varepsilon$  is some positive number. The neighborhood is also called a *ball* with radius  $\varepsilon$  and center  $\mathbf{x}$ .

In the plane  $\mathbb{R}^2$ , a neighborhood of  $\mathbf{x} = [x_1, x_2]^\top$  consists of all the points inside a disk centered at  $\mathbf{x}$ . In  $\mathbb{R}^3$ , a neighborhood of  $\mathbf{x} = [x_1, x_2, x_3]^\top$  consists of all the points inside a sphere centered at  $\mathbf{x}$  (see Figure 4.7).



**Figure 4.7** Examples of neighborhoods of a point in  $\mathbb{R}^2$  and  $\mathbb{R}^3$ .

A point  $x \in S$  is said to be an *interior point* of the set  $S$  if the set  $S$  contains some neighborhood of  $x$ ; that is, if all points within some neighborhood of  $x$  are also in  $S$  (see Figure 4.8). The set of all the interior points of  $S$  is called the *interior* of  $S$ .

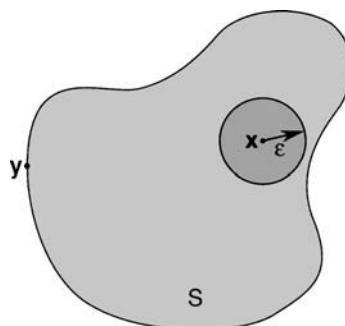
A point  $x$  is said to be a *boundary point* of the set  $S$  if every neighborhood of  $x$  contains a point in  $S$  and a point not in  $S$  (see Figure 4.8). Note that a boundary point of  $S$  may or may not be an element of  $S$ . The set of all boundary points of  $S$  is called the *boundary* of  $S$ .

A set  $S$  is said to be *open* if it contains a neighborhood of each of its points; that is, if each of its points is an interior point, or equivalently, if  $S$  contains no boundary points.

A set  $S$  is said to be *closed* if it contains its boundary (see Figure 4.9). We can show that a set is closed if and only if its complement is open.

A set that is contained in a ball of finite radius is said to be *bounded*. A set is *compact* if it is both closed and bounded. Compact sets are important in optimization problems for the following reason.

**Theorem 4.2 Theorem of Weierstrass.** Let  $f : \Omega \rightarrow \mathbb{R}$  be a continuous function, where  $\Omega \subset \mathbb{R}^n$  is a compact set. Then, there exists a point  $x_0 \in \Omega$



**Figure 4.8**  $x$  is an interior point;  $y$  is a boundary point.

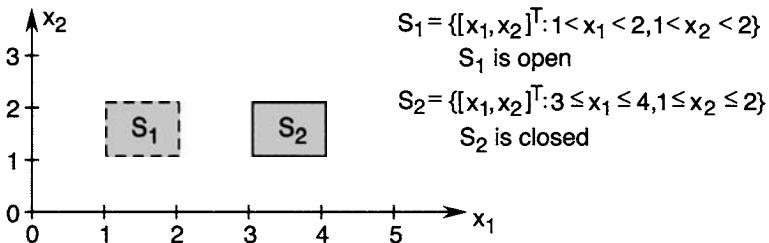


Figure 4.9 Open and closed sets.

such that  $f(\mathbf{x}_0) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \Omega$ . In other words,  $f$  achieves its minimum on  $\Omega$ .  $\square$

*Proof.* See [112, p. 89] or [2, p. 154].  $\blacksquare$

## 4.5 Polytopes and Polyhedra

Let  $\Theta$  be a convex set, and suppose that  $\mathbf{y}$  is a boundary point of  $\Theta$ . A hyperplane passing through  $\mathbf{y}$  is called a *hyperplane of support* (or *supporting hyperplane*) of the set  $\Theta$  if the entire set  $\Theta$  lies completely in one of the two half-spaces into which this hyperplane divides the space  $\mathbb{R}^n$ .

Recall that by Theorem 4.1, the intersection of any number of convex sets is convex. In what follows we are concerned with the intersection of a finite number of half-spaces. Because every half-space  $H_+$  or  $H_-$  is convex in  $\mathbb{R}^n$ , the intersection of any number of half-spaces is a convex set.

A set that can be expressed as the intersection of a finite number of half-spaces is called a *convex polytope* (see Figure 4.10).

A nonempty bounded polytope is called a *polyhedron* (see Figure 4.11).

For every convex polyhedron  $\Theta \subset \mathbb{R}^n$ , there exists a nonnegative integer  $k \leq n$  such that  $\Theta$  is contained in a linear variety of dimension  $k$ , but is not

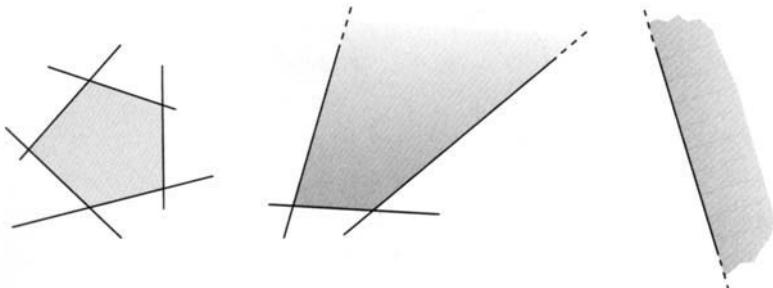
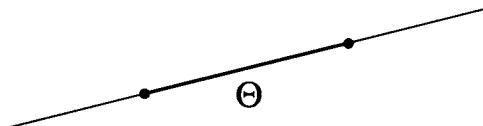


Figure 4.10 Polytopes.



**Figure 4.11** One-dimensional polyhedron.

entirely contained in any  $(k - 1)$ -dimensional linear variety of  $\mathbb{R}^n$ . Furthermore, there exists only one  $k$ -dimensional linear variety containing  $\Theta$ , called the *carrier* of the polyhedron  $\Theta$ , and  $k$  is called the dimension of  $\Theta$ . For example, a zero-dimensional polyhedron is a point of  $\mathbb{R}^n$ , and its carrier is itself. A one-dimensional polyhedron is a segment, and its carrier is the straight line on which it lies. The boundary of any  $k$ -dimensional polyhedron,  $k > 0$ , consists of a finite number of  $(k - 1)$ -dimensional polyhedra. For example, the boundary of a one-dimensional polyhedron consists of two points that are the endpoints of the segment.

The  $(k - 1)$ -dimensional polyhedra forming the boundary of a  $k$ -dimensional polyhedron are called the *faces* of the polyhedron. Each of these faces has, in turn,  $(k - 2)$ -dimensional faces. We also consider each of these  $(k - 2)$ -dimensional faces to be faces of the original  $k$ -dimensional polyhedron. Thus, every  $k$ -dimensional polyhedron has faces of dimensions  $k - 1, k - 2, \dots, 1, 0$ . A zero-dimensional face of a polyhedron is called a *vertex*, and a one-dimensional face is called an *edge*.

## EXERCISES

**4.1** Show that a set  $S \subset \mathbb{R}^n$  is a linear variety if and only if for all  $\mathbf{x}, \mathbf{y} \in S$  and  $\alpha \in \mathbb{R}$ , we have  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in S$ .

**4.2** Show that the set  $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq r\}$  is convex, where  $r > 0$  is a given real number and  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}}$  is the Euclidean norm of  $\mathbf{x} \in \mathbb{R}^n$ .

**4.3** Show that for any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and vector  $\mathbf{b} \in \mathbb{R}^m$ , the set (linear variety)  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} = \mathbf{b}\}$  is convex.

**4.4** Show that the set  $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}\}$  is convex (where  $\mathbf{x} \geq \mathbf{0}$  means that every component of  $\mathbf{x}$  is nonnegative).



# CHAPTER 5

---

## ELEMENTS OF CALCULUS

---

### 5.1 Sequences and Limits

A *sequence of real numbers* is a function whose domain is the set of natural numbers  $1, 2, \dots, k, \dots$  and whose range is contained in  $\mathbb{R}$ . Thus, a sequence of real numbers can be viewed as a set of numbers  $\{x_1, x_2, \dots, x_k, \dots\}$ , which is often also denoted as  $\{x_k\}$  (or sometimes as  $\{x_k\}_{k=1}^{\infty}$ , to indicate explicitly the range of values that  $k$  can take).

A sequence  $\{x_k\}$  is *increasing* if  $x_1 < x_2 < \dots < x_k \dots$ ; that is,  $x_k < x_{k+1}$  for all  $k$ . If  $x_k \leq x_{k+1}$ , then we say that the sequence is *nondecreasing*. Similarly, we can define *decreasing* and *nonincreasing sequences*. Nonincreasing or nondecreasing sequences are called *monotone sequences*.

A number  $x^* \in \mathbb{R}$  is called the *limit* of the sequence  $\{x_k\}$  if for any positive  $\varepsilon$  there is a number  $K$  (which may depend on  $\varepsilon$ ) such that for all  $k > K$ ,  $|x_k - x^*| < \varepsilon$ ; that is,  $x_k$  lies between  $x^* - \varepsilon$  and  $x^* + \varepsilon$  for all  $k > K$ . In this case we write

$$x^* = \lim_{k \rightarrow \infty} x_k$$

or

$$x_k \rightarrow x^*.$$

A sequence that has a limit is called a *convergent sequence*.

The notion of a sequence can be extended to sequences with elements in  $\mathbb{R}^n$ . Specifically, a sequence in  $\mathbb{R}^n$  is a function whose domain is the set of natural numbers  $1, 2, \dots, k, \dots$  and whose range is contained in  $\mathbb{R}^n$ . We use the notation  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots\}$  or  $\{\mathbf{x}^{(k)}\}$  for sequences in  $\mathbb{R}^n$ . For limits of sequences in  $\mathbb{R}^n$ , we need to replace absolute values with vector norms. In other words,  $\mathbf{x}^*$  is the limit of  $\{\mathbf{x}^{(k)}\}$  if for any positive  $\varepsilon$  there is a number  $K$  (which may depend on  $\varepsilon$ ) such that for all  $k > K$ ,  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \varepsilon$ . As before, if a sequence  $\{\mathbf{x}^{(k)}\}$  is convergent, we write  $\mathbf{x}^* = \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$  or  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ .

**Theorem 5.1** *A convergent sequence has only one limit.* □

*Proof.* We prove this result by contradiction. Suppose that a sequence  $\{\mathbf{x}^{(k)}\}$  has two different limits, say  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Then, we have  $\|\mathbf{x}_1 - \mathbf{x}_2\| > 0$ . Let

$$\varepsilon = \frac{1}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

From the definition of a limit, there exist  $K_1$  and  $K_2$  such that for  $k > K_1$  we have  $\|\mathbf{x}^{(k)} - \mathbf{x}_1\| < \varepsilon$ , and for  $k > K_2$  we have  $\|\mathbf{x}^{(k)} - \mathbf{x}_2\| < \varepsilon$ . Let  $K = \max\{K_1, K_2\}$ . Then, if  $k > K$ , we have  $\|\mathbf{x}^{(k)} - \mathbf{x}_1\| < \varepsilon$  and  $\|\mathbf{x}^{(k)} - \mathbf{x}_2\| < \varepsilon$ . Adding  $\|\mathbf{x}^{(k)} - \mathbf{x}_1\| < \varepsilon$  and  $\|\mathbf{x}^{(k)} - \mathbf{x}_2\| < \varepsilon$  yields

$$\|\mathbf{x}^{(k)} - \mathbf{x}_1\| + \|\mathbf{x}^{(k)} - \mathbf{x}_2\| < 2\varepsilon.$$

Applying the triangle inequality gives

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}_2\| &= \|\mathbf{x}^{(k)} - \mathbf{x}_1 - (\mathbf{x}^{(k)} - \mathbf{x}_2)\| \\ &= \|(\mathbf{x}^{(k)} - \mathbf{x}_1) - (\mathbf{x}^{(k)} - \mathbf{x}_2)\| \\ &\leq \|\mathbf{x}^{(k)} - \mathbf{x}_1\| + \|\mathbf{x}^{(k)} - \mathbf{x}_2\|. \end{aligned}$$

Therefore,

$$\|\mathbf{x}_1 - \mathbf{x}_2\| = \|\mathbf{x}_1 - \mathbf{x}_2\| < 2\varepsilon.$$

However, this contradicts the assumption that  $\|\mathbf{x}_1 - \mathbf{x}_2\| = 2\varepsilon$ , which completes the proof. ■

A sequence  $\{\mathbf{x}^{(k)}\}$  in  $\mathbb{R}^n$  is *bounded* if there exists a number  $B \geq 0$  such that  $\|\mathbf{x}^{(k)}\| \leq B$  for all  $k = 1, 2, \dots$ .

**Theorem 5.2** *Every convergent sequence is bounded.* □

*Proof.* Let  $\{\mathbf{x}^{(k)}\}$  be a convergent sequence with limit  $\mathbf{x}^*$ . Choose  $\varepsilon = 1$ . Then, by definition of the limit, there exists a natural number  $K$  such that for all  $k > K$ ,

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < 1.$$

By the result of Exercise 2.9, we get

$$\|\mathbf{x}^{(k)}\| - \|\mathbf{x}^*\| \leq \|\mathbf{x}^{(k)} - \mathbf{x}^*\| < 1 \quad \text{for all } k > K.$$

Therefore,

$$\|\mathbf{x}^{(k)}\| < \|\mathbf{x}^*\| + 1 \quad \text{for all } k > K.$$

Letting

$$B = \max \left\{ \|\mathbf{x}^{(1)}\|, \|\mathbf{x}^{(2)}\|, \dots, \|\mathbf{x}^{(K)}\|, \|\mathbf{x}^*\| + 1 \right\},$$

we have

$$B \geq \|\mathbf{x}^{(k)}\| \quad \text{for all } k,$$

which means that the sequence  $\{\mathbf{x}^{(k)}\}$  is bounded. ■

For a sequence  $\{x_k\}$  in  $\mathbb{R}$ , a number  $B$  is called an *upper bound* if  $x_k \leq B$  for all  $k = 1, 2, \dots$ . In this case, we say that  $\{x_k\}$  is *bounded above*. Similarly,  $B$  is called a *lower bound* if  $x_k \geq B$  for all  $k = 1, 2, \dots$ . In this case, we say that  $\{x_k\}$  is *bounded below*. Clearly, a sequence is bounded if it is both bounded above and bounded below.

Any sequence  $\{x_k\}$  in  $\mathbb{R}$  that has an upper bound has a *least upper bound* (also called the *supremum*), which is the smallest number  $B$  that is an upper bound of  $\{x_k\}$ . Similarly, any sequence  $\{x_k\}$  in  $\mathbb{R}$  that has a lower bound has a *greatest lower bound* (also called the *infimum*). If  $B$  is the least upper bound of the sequence  $\{x_k\}$ , then  $x_k \leq B$  for all  $k$ , and for any  $\varepsilon > 0$ , there exists a number  $K$  such that  $x_K > B - \varepsilon$ . An analogous statement applies to the greatest lower bound: If  $B$  is the greatest lower bound of  $\{x_k\}$ , then  $x_k \geq B$  for all  $k$ , and for any  $\varepsilon > 0$ , there exists a number  $K$  such that  $x_K < B + \varepsilon$ .

**Theorem 5.3** *Every monotone bounded sequence in  $\mathbb{R}$  is convergent.* □

*Proof.* We prove the theorem for nondecreasing sequences. The proof for nonincreasing sequences is analogous.

Let  $\{x_k\}$  be a bounded nondecreasing sequence in  $\mathbb{R}$  and  $x^*$  the least upper bound. Fix a number  $\varepsilon > 0$ . Then, there exists a number  $K$  such that  $x_K > x^* - \varepsilon$ . Because  $\{x_k\}$  is nondecreasing, for any  $k \geq K$ ,

$$x_k \geq x_K > x^* - \varepsilon.$$

Also, because  $x^*$  is an upper bound of  $\{x_k\}$ , we have

$$x_k \leq x^* < x^* + \varepsilon.$$

Therefore, for any  $k \geq K$ ,

$$|x_k - x^*| < \varepsilon,$$

which means that  $x_k \rightarrow x^*$ . ■

Suppose that we are given a sequence  $\{\mathbf{x}^{(k)}\}$  and an increasing sequence of natural numbers  $\{m_k\}$ . The sequence

$$\{\mathbf{x}^{(m_k)}\} = \{\mathbf{x}^{(m_1)}, \mathbf{x}^{(m_2)}, \dots\}$$

is called a *subsequence* of the sequence  $\{\mathbf{x}^{(k)}\}$ . A subsequence of a given sequence can thus be obtained by neglecting some elements of the given sequence.

**Theorem 5.4** *Consider a convergent sequence  $\{\mathbf{x}^{(k)}\}$  with limit  $\mathbf{x}^*$ . Then, any subsequence of  $\{\mathbf{x}^{(k)}\}$  also converges to  $\mathbf{x}^*$ .*  $\square$

*Proof.* Let  $\{\mathbf{x}^{(m_k)}\}$  be a subsequence of  $\{\mathbf{x}^{(k)}\}$ , where  $\{m_k\}$  is an increasing sequence of natural numbers. Observe that  $m_k \geq k$  for all  $k = 1, 2, \dots$ . To show this, first note that  $m_1 \geq 1$  because  $m_1$  is a natural number. Next, we proceed by induction by assuming that  $m_k \geq k$ . Then, we have  $m_{k+1} > m_k \geq k$ , which implies that  $m_{k+1} \geq k + 1$ . Therefore, we have shown that  $m_k \geq k$  for all  $k = 1, 2, \dots$

Let  $\varepsilon > 0$  be given. Then, by definition of the limit, there exists  $K$  such that  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| < \varepsilon$  for any  $k > K$ . Because  $m_k \geq k$ , we also have  $\|\mathbf{x}^{(m_k)} - \mathbf{x}^*\| < \varepsilon$  for any  $k > K$ . This means that

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(m_k)} = \mathbf{x}^*. \quad \blacksquare$$

It turns out that any bounded sequence contains a convergent subsequence. This result is called the *Bolzano-Weierstrass theorem* (see [2, p. 70]).

Consider a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a point  $\mathbf{x}_0 \in \mathbb{R}^n$ . Suppose that there exists  $\mathbf{f}^*$  such that for any convergent sequence  $\{\mathbf{x}^{(k)}\}$  with limit  $\mathbf{x}_0$ , we have

$$\lim_{k \rightarrow \infty} \mathbf{f}(\mathbf{x}^{(k)}) = \mathbf{f}^*.$$

Then, we use the notation

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_0} \mathbf{f}(\mathbf{x})$$

to represent the limit  $\mathbf{f}^*$ .

It turns out that  $\mathbf{f}$  is continuous at  $\mathbf{x}_0$  if and only if for any convergent sequence  $\{\mathbf{x}^{(k)}\}$  with limit  $\mathbf{x}_0$ , we have

$$\lim_{k \rightarrow \infty} \mathbf{f}(\mathbf{x}^{(k)}) = \mathbf{f}\left(\lim_{k \rightarrow \infty} \mathbf{x}^{(k)}\right) = \mathbf{f}(\mathbf{x}_0)$$

(see [2, p. 137]). Therefore, using the notation introduced above, the function  $\mathbf{f}$  is continuous at  $\mathbf{x}_0$  if and only if

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_0} \mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_0).$$

We end this section with some results involving sequences and limits of matrices. These results are useful in the analysis of algorithms (e.g., Newton's algorithm in Chapter 9).

We say that a sequence  $\{\mathbf{A}_k\}$  of  $m \times n$  matrices converges to the  $m \times n$  matrix  $\mathbf{A}$  if

$$\lim_{k \rightarrow \infty} \|\mathbf{A} - \mathbf{A}_k\| = 0.$$

**Lemma 5.1** *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then,  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{O}$  if and only if the eigenvalues of  $\mathbf{A}$  satisfy  $|\lambda_i(\mathbf{A})| < 1$ ,  $i = 1, \dots, n$ .*  $\square$

*Proof.* To prove this theorem, we use the *Jordan form* (see, e.g., [47]). Specifically, it is well known that any square matrix is similar to the Jordan form: There exists a nonsingular  $\mathbf{T}$  such that

$$\mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \text{diag} [\mathbf{J}_{m_1}(\lambda_1), \dots, \mathbf{J}_{m_s}(\lambda_1), \mathbf{J}_{n_1}(\lambda_2), \dots, \mathbf{J}_{t_\nu}(\lambda_q)] \triangleq \mathbf{J},$$

where  $\mathbf{J}_r(\lambda)$  is the  $r \times r$  matrix:

$$\mathbf{J}_r(\lambda) = \begin{bmatrix} \lambda & 1 & & 0 \\ & \lambda & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda \end{bmatrix}.$$

The  $\lambda_1, \dots, \lambda_q$  above are distinct eigenvalues of  $\mathbf{A}$ , the multiplicity of  $\lambda_1$  is  $m_1 + \dots + m_s$ , and so on.

We may rewrite the above as  $\mathbf{A} = \mathbf{T}^{-1}\mathbf{J}\mathbf{T}$ . To complete the proof, observe that

$$(\mathbf{J}_r(\lambda))^k = \begin{bmatrix} \lambda^k & \binom{k}{k-1} \lambda^{k-1} & \dots & \binom{k}{k-r+1} \lambda^{k-r+1} \\ 0 & \lambda^k & \dots & \binom{k}{k-r+2} \lambda^{k-r+2} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda^k \end{bmatrix},$$

where

$$\binom{k}{i} = \frac{k!}{i!(k-i)!}.$$

Furthermore,

$$\mathbf{A}^k = \mathbf{T}^{-1}\mathbf{J}^k\mathbf{T}.$$

Hence,

$$\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{T}^{-1} \left( \lim_{k \rightarrow \infty} \mathbf{J}^k \right) \mathbf{T} = \mathbf{O}$$

if and only if  $|\lambda_i| < 1$ ,  $i = 1, \dots, n$ . ■

**Lemma 5.2** *The series of  $n \times n$  matrices*

$$\mathbf{I}_n + \mathbf{A} + \mathbf{A}^2 + \cdots + \mathbf{A}^k + \cdots$$

*converges if and only if  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{O}$ . In this case the sum of the series equals  $(\mathbf{I}_n - \mathbf{A})^{-1}$ .* □

*Proof.* The necessity of the condition is obvious.

To prove the sufficiency, suppose that  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{O}$ . By Lemma 5.1 we deduce that  $|\lambda_i(\mathbf{A})| < 1$ ,  $i = 1, \dots, n$ . This implies that  $\det(\mathbf{I}_n - \mathbf{A}) \neq 0$ , and hence  $(\mathbf{I}_n - \mathbf{A})^{-1}$  exists. Consider now the following relation:

$$(\mathbf{I}_n + \mathbf{A} + \mathbf{A}^2 + \cdots + \mathbf{A}^k)(\mathbf{I}_n - \mathbf{A}) = \mathbf{I}_n - \mathbf{A}^{k+1}.$$

Postmultiplying the equation above by  $(\mathbf{I}_n - \mathbf{A})^{-1}$  yields

$$\mathbf{I}_n + \mathbf{A} + \mathbf{A}^2 + \cdots + \mathbf{A}^k = (\mathbf{I}_n - \mathbf{A})^{-1} - \mathbf{A}^{k+1}(\mathbf{I}_n - \mathbf{A})^{-1}.$$

Hence,

$$\lim_{k \rightarrow \infty} \sum_{j=0}^k \mathbf{A}^j = (\mathbf{I}_n - \mathbf{A})^{-1},$$

because  $\lim_{k \rightarrow \infty} \mathbf{A}^{k+1} = \mathbf{O}$ . Thus,

$$\sum_{j=0}^{\infty} \mathbf{A}^j = (\mathbf{I}_n - \mathbf{A})^{-1},$$

which completes the proof. ■

A matrix-valued function  $\mathbf{A} : \mathbb{R}^r \rightarrow \mathbb{R}^{n \times n}$  is continuous at a point  $\xi_0 \in \mathbb{R}^r$  if

$$\lim_{\|\xi - \xi_0\| \rightarrow 0} \|\mathbf{A}(\xi) - \mathbf{A}(\xi_0)\| = 0.$$

**Lemma 5.3** *Let  $\mathbf{A} : \mathbb{R}^r \rightarrow \mathbb{R}^{n \times n}$  be an  $n \times n$  matrix-valued function that is continuous at  $\xi_0$ . If  $\mathbf{A}(\xi_0)^{-1}$  exists, then  $\mathbf{A}(\xi)^{-1}$  exists for  $\xi$  sufficiently close to  $\xi_0$  and  $\mathbf{A}(\cdot)^{-1}$  is continuous at  $\xi_0$ .* □

*Proof.* We follow [114]. We first prove the existence of  $\mathbf{A}(\xi)^{-1}$  for all  $\xi$  sufficiently close to  $\xi_0$ . We have

$$\mathbf{A}(\xi) = \mathbf{A}(\xi_0) - \mathbf{A}(\xi_0) + \mathbf{A}(\xi) = \mathbf{A}(\xi_0)(\mathbf{I}_n - \mathbf{K}(\xi)),$$

where

$$\mathbf{K}(\xi) = \mathbf{A}(\xi_0)^{-1}(\mathbf{A}(\xi_0) - \mathbf{A}(\xi)).$$

Thus,

$$\|\mathbf{K}(\boldsymbol{\xi})\| \leq \|\mathbf{A}(\boldsymbol{\xi}_0)^{-1}\| \|\mathbf{A}(\boldsymbol{\xi}_0) - \mathbf{A}(\boldsymbol{\xi})\|$$

and

$$\lim_{\|\boldsymbol{\xi} - \boldsymbol{\xi}_0\| \rightarrow 0} \|\mathbf{K}(\boldsymbol{\xi})\| = 0.$$

Because  $\mathbf{A}$  is continuous at  $\boldsymbol{\xi}_0$ , for all  $\boldsymbol{\xi}$  close enough to  $\boldsymbol{\xi}_0$ , we have

$$\|\mathbf{A}(\boldsymbol{\xi}_0) - \mathbf{A}(\boldsymbol{\xi})\| \leq \frac{\theta}{\|\mathbf{A}(\boldsymbol{\xi}_0)^{-1}\|},$$

where  $\theta \in (0, 1)$ . Then,

$$\|\mathbf{K}(\boldsymbol{\xi})\| \leq \theta < 1$$

and

$$(\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi}))^{-1}$$

exists. But then

$$\mathbf{A}(\boldsymbol{\xi})^{-1} = (\mathbf{A}(\boldsymbol{\xi}_0)(\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi})))^{-1} = (\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi}))^{-1} \mathbf{A}(\boldsymbol{\xi}_0)^{-1},$$

which means that  $\mathbf{A}(\boldsymbol{\xi})^{-1}$  exists for  $\boldsymbol{\xi}$  sufficiently close to  $\boldsymbol{\xi}_0$ .

To prove the continuity of  $\mathbf{A}(\cdot)^{-1}$  note that

$$\begin{aligned} \|\mathbf{A}(\boldsymbol{\xi}_0)^{-1} - \mathbf{A}(\boldsymbol{\xi})^{-1}\| &= \|\mathbf{A}(\boldsymbol{\xi})^{-1} - \mathbf{A}(\boldsymbol{\xi}_0)^{-1}\| \\ &= \|((\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi}))^{-1} - \mathbf{I}_n)\mathbf{A}(\boldsymbol{\xi}_0)^{-1}\|. \end{aligned}$$

However, since  $\|\mathbf{K}(\boldsymbol{\xi})\| < 1$ , it follows from Lemma 5.2 that

$$(\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi}))^{-1} - \mathbf{I}_n = \mathbf{K}(\boldsymbol{\xi}) + \mathbf{K}^2(\boldsymbol{\xi}) + \cdots = \mathbf{K}(\boldsymbol{\xi})(\mathbf{I}_n + \mathbf{K}(\boldsymbol{\xi}) + \cdots).$$

Hence,

$$\begin{aligned} \|(\mathbf{I}_n - \mathbf{K}(\boldsymbol{\xi}))^{-1} - \mathbf{I}_n\| &\leq \|\mathbf{K}(\boldsymbol{\xi})\|(1 + \|\mathbf{K}(\boldsymbol{\xi})\| + \|\mathbf{K}(\boldsymbol{\xi})\|^2 + \cdots) \\ &= \frac{\|\mathbf{K}(\boldsymbol{\xi})\|}{1 - \|\mathbf{K}(\boldsymbol{\xi})\|}, \end{aligned}$$

when  $\|\mathbf{K}(\boldsymbol{\xi})\| < 1$ . Therefore,

$$\|\mathbf{A}(\boldsymbol{\xi})^{-1} - \mathbf{A}(\boldsymbol{\xi}_0)^{-1}\| \leq \frac{\|\mathbf{K}(\boldsymbol{\xi})\|}{1 - \|\mathbf{K}(\boldsymbol{\xi})\|} \|\mathbf{A}(\boldsymbol{\xi}_0)^{-1}\|.$$

Because

$$\lim_{\|\boldsymbol{\xi} - \boldsymbol{\xi}_0\| \rightarrow 0} \|\mathbf{K}(\boldsymbol{\xi})\| = 0,$$

we obtain

$$\lim_{\|\boldsymbol{\xi} - \boldsymbol{\xi}_0\| \rightarrow 0} \|\mathbf{A}(\boldsymbol{\xi})^{-1} - \mathbf{A}(\boldsymbol{\xi}_0)^{-1}\| = 0,$$

which completes the proof. ■

## 5.2 Differentiability

Differential calculus is based on the idea of approximating an arbitrary function by an *affine function*. A function  $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is *affine* if there exists a *linear function*  $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a vector  $\mathbf{y} \in \mathbb{R}^m$  such that

$$\mathcal{A}(\mathbf{x}) = \mathcal{L}(\mathbf{x}) + \mathbf{y}$$

for every  $\mathbf{x} \in \mathbb{R}^n$ . Consider a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a point  $\mathbf{x}_0 \in \mathbb{R}^n$ . We wish to find an affine function  $\mathcal{A}$  that approximates  $\mathbf{f}$  near the point  $\mathbf{x}_0$ . First, it is natural to impose the condition

$$\mathcal{A}(\mathbf{x}_0) = \mathbf{f}(\mathbf{x}_0).$$

Because  $\mathcal{A}(\mathbf{x}) = \mathcal{L}(\mathbf{x}) + \mathbf{y}$ , we obtain  $\mathbf{y} = \mathbf{f}(\mathbf{x}_0) - \mathcal{L}(\mathbf{x}_0)$ . By the linearity of  $\mathcal{L}$ ,

$$\mathcal{L}(\mathbf{x}) + \mathbf{y} = \mathcal{L}(\mathbf{x}) - \mathcal{L}(\mathbf{x}_0) + \mathbf{f}(\mathbf{x}_0) = \mathcal{L}(\mathbf{x} - \mathbf{x}_0) + \mathbf{f}(\mathbf{x}_0).$$

Hence, we may write

$$\mathcal{A}(\mathbf{x}) = \mathcal{L}(\mathbf{x} - \mathbf{x}_0) + \mathbf{f}(\mathbf{x}_0).$$

Next, we require that  $\mathcal{A}(\mathbf{x})$  approaches  $\mathbf{f}(\mathbf{x})$  faster than  $\mathbf{x}$  approaches  $\mathbf{x}_0$ ; that is,

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_0, \mathbf{x} \in \Omega} \frac{\|\mathbf{f}(\mathbf{x}) - \mathcal{A}(\mathbf{x})\|}{\|\mathbf{x} - \mathbf{x}_0\|} = 0.$$

The conditions above on  $\mathcal{A}$  ensure that  $\mathcal{A}$  approximates  $\mathbf{f}$  near  $\mathbf{x}_0$  in the sense that the error in the approximation at a given point is “small” compared with the distance of the point from  $\mathbf{x}_0$ .

In summary, a function  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^m$ ,  $\Omega \subset \mathbb{R}^n$ , is said to be *differentiable* at  $\mathbf{x}_0 \in \Omega$  if there is an affine function that approximates  $\mathbf{f}$  near  $\mathbf{x}_0$ ; that is, there exists a linear function  $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_0, \mathbf{x} \in \Omega} \frac{\|\mathbf{f}(\mathbf{x}) - (\mathcal{L}(\mathbf{x} - \mathbf{x}_0) + \mathbf{f}(\mathbf{x}_0))\|}{\|\mathbf{x} - \mathbf{x}_0\|} = 0.$$

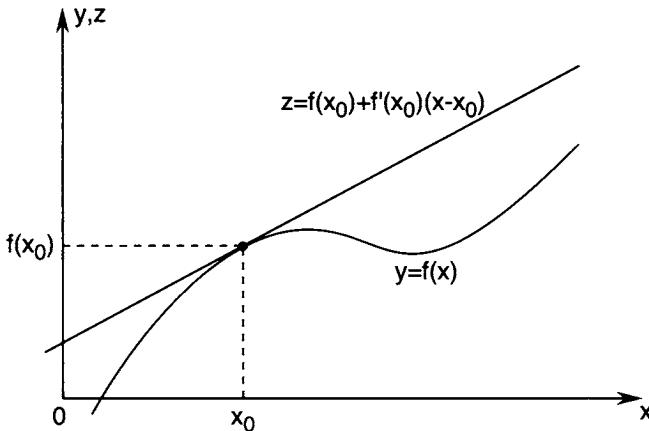
The linear function  $\mathcal{L}$  above is determined uniquely by  $\mathbf{f}$  and  $\mathbf{x}_0$  and is called the *derivative* of  $\mathbf{f}$  at  $\mathbf{x}_0$ . The function  $\mathbf{f}$  is said to be *differentiable* on  $\Omega$  if  $\mathbf{f}$  is differentiable at every point of its domain  $\Omega$ .

In  $\mathbb{R}$ , an affine function has the form  $ax + b$ , with  $a, b \in \mathbb{R}$ . Hence, a real-valued function  $f(x)$  of a real variable  $x$  that is differentiable at  $x_0$  can be approximated near  $x_0$  by a function

$$\mathcal{A}(x) = ax + b.$$

Because  $f(x_0) = \mathcal{A}(x_0) = ax_0 + b$ , we obtain

$$\mathcal{A}(x) = ax + b = a(x - x_0) + f(x_0).$$



**Figure 5.1** Illustration of the notion of the derivative.

The linear part of  $\mathcal{A}(x)$ , denoted earlier by  $\mathcal{L}(x)$ , is in this case just  $ax$ . The norm of a real number is its absolute value, so by the definition of differentiability we have

$$\lim_{x \rightarrow x_0} \frac{|f(x) - (a(x - x_0) + f(x_0))|}{|x - x_0|} = 0,$$

which is equivalent to

$$\lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = a.$$

The number  $a$  is commonly denoted  $f'(x_0)$  and is called the derivative of  $f$  at  $x_0$ . The affine function  $\mathcal{A}$  is therefore given by

$$\mathcal{A}(x) = f(x_0) + f'(x_0)(x - x_0).$$

This affine function is tangent to  $f$  at  $x_0$  (see Figure 5.1).

### 5.3 The Derivative Matrix

Any linear transformation from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ , and in particular the derivative  $\mathcal{L}$  of  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , can be represented by an  $m \times n$  matrix. To find the matrix representation  $\mathbf{L}$  of the derivative  $\mathcal{L}$  of a differentiable function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , we use the natural basis  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  for  $\mathbb{R}^n$ . Consider the vectors

$$\mathbf{x}_j = \mathbf{x}_0 + t\mathbf{e}_j, \quad j = 1, \dots, n.$$

By the definition of the derivative, we have

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x}_j) - (t\mathbf{L}\mathbf{e}_j + \mathbf{f}(\mathbf{x}_0))}{t} = \mathbf{0}$$

for  $j = 1, \dots, n$ . This means that

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x}_j) - \mathbf{f}(\mathbf{x}_0)}{t} = \mathbf{L}\mathbf{e}_j$$

for  $j = 1, \dots, n$ . But  $\mathbf{L}\mathbf{e}_j$  is the  $j$ th column of the matrix  $\mathbf{L}$ . On the other hand, the vector  $\mathbf{x}_j$  differs from  $\mathbf{x}_0$  only in the  $j$ th coordinate, and in that coordinate the difference is just the number  $t$ . Therefore, the left side of the preceding equation is the partial derivative

$$\frac{\partial \mathbf{f}}{\partial x_j}(\mathbf{x}_0).$$

Because vector limits are computed by taking the limit of each coordinate function, it follows that if

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix},$$

then

$$\frac{\partial \mathbf{f}}{\partial x_j}(\mathbf{x}_0) = \begin{bmatrix} \frac{\partial f_1}{\partial x_j}(\mathbf{x}_0) \\ \vdots \\ \frac{\partial f_m}{\partial x_j}(\mathbf{x}_0) \end{bmatrix},$$

and the matrix  $\mathbf{L}$  has the form

$$\left[ \frac{\partial \mathbf{f}}{\partial x_1}(\mathbf{x}_0) \ \cdots \ \frac{\partial \mathbf{f}}{\partial x_n}(\mathbf{x}_0) \right] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}_0) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}_0) \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}_0) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}_0) \end{bmatrix}.$$

The matrix  $\mathbf{L}$  is called the *Jacobian matrix*, or *derivative matrix*, of  $\mathbf{f}$  at  $\mathbf{x}_0$ , and is denoted  $D\mathbf{f}(\mathbf{x}_0)$ . For convenience, we often refer to  $D\mathbf{f}(\mathbf{x}_0)$  simply as the derivative of  $\mathbf{f}$  at  $\mathbf{x}_0$ . We summarize the foregoing discussion in the following theorem.

**Theorem 5.5** *If a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable at  $\mathbf{x}_0$ , then the derivative of  $\mathbf{f}$  at  $\mathbf{x}_0$  is determined uniquely and is represented by the  $m \times n$  derivative matrix  $D\mathbf{f}(\mathbf{x}_0)$ . The best affine approximation to  $\mathbf{f}$  near  $\mathbf{x}_0$  is then given by*

$$\mathcal{A}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_0) + D\mathbf{f}(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0),$$

*in the sense that*

$$\mathbf{f}(\mathbf{x}) = \mathcal{A}(\mathbf{x}) + \mathbf{r}(\mathbf{x})$$

and  $\lim_{\mathbf{x} \rightarrow \mathbf{x}_0} \|\mathbf{r}(\mathbf{x})\|/\|\mathbf{x} - \mathbf{x}_0\| = 0$ . The columns of the derivative matrix  $D\mathbf{f}(\mathbf{x}_0)$  are vector partial derivatives. The vector

$$\frac{\partial \mathbf{f}}{\partial x_j}(\mathbf{x}_0)$$

is a tangent vector at  $\mathbf{x}_0$  to the curve  $\mathbf{f}$  obtained by varying only the  $j$ th coordinate of  $\mathbf{x}$ .  $\square$

If  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable, then the function  $\nabla f$  defined by

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial f}{\partial x_n}(\mathbf{x}) \end{bmatrix} = Df(\mathbf{x})^\top$$

is called the *gradient* of  $f$ . The gradient is a function from  $\mathbb{R}^n$  to  $\mathbb{R}^n$ , and can be pictured as a *vector field*, by drawing the arrow representing  $\nabla f(\mathbf{x})$  so that its tail starts at  $\mathbf{x}$ .

Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , if  $\nabla f$  is differentiable, we say that  $f$  is *twice differentiable*, and we write the derivative of  $\nabla f$  as

$$D^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

(The notation  $\frac{\partial^2 f}{\partial x_i \partial x_j}$  represents taking the partial derivative of  $f$  with respect to  $x_j$  first, then with respect to  $x_i$ .) The matrix  $D^2 f(\mathbf{x})$  is called the *Hessian matrix* of  $f$  at  $\mathbf{x}$ , and is often also denoted  $\mathbf{F}(\mathbf{x})$ .

A function  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^m$ ,  $\Omega \subset \mathbb{R}^n$ , is said to be *continuously differentiable* on  $\Omega$  if it is differentiable (on  $\Omega$ ), and  $D\mathbf{f} : \Omega \rightarrow \mathbb{R}^{m \times n}$  is continuous; that is, the components of  $\mathbf{f}$  have continuous partial derivatives. In this case, we write  $\mathbf{f} \in C^1$ . If the components of  $\mathbf{f}$  have continuous partial derivatives of order  $p$ , then we write  $\mathbf{f} \in C^p$ .

Note that the Hessian matrix of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  at  $\mathbf{x}$  is symmetric if  $f$  is twice continuously differentiable at  $\mathbf{x}$ . This is a well-known result from calculus called *Clairaut's theorem* or *Schwarz's theorem*. However, if the second partial derivatives of  $f$  are not continuous, then there is no guarantee that the Hessian is symmetric, as shown in the following well-known example.

**Example 5.1** Consider the function

$$f(\mathbf{x}) = \begin{cases} x_1 x_2 (x_1^2 - x_2^2)/(x_1^2 + x_2^2) & \text{if } \mathbf{x} \neq \mathbf{0} \\ 0 & \text{if } \mathbf{x} = \mathbf{0}. \end{cases}$$

Let us compute its Hessian at the point  $\mathbf{0} = [0, 0]^\top$ . We have

$$\mathbf{F} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}.$$

We now proceed with computing the components of the Hessian and evaluating them at the point  $[0, 0]^\top$  one by one. We start with

$$\frac{\partial^2 f}{\partial x_1^2} = \frac{\partial}{\partial x_1} \left( \frac{\partial f}{\partial x_1} \right),$$

where

$$\begin{aligned} \frac{\partial f}{\partial x_1}(\mathbf{x}) \\ = \begin{cases} x_2(x_1^4 - x_2^4 + 4x_1^2x_2^2)/(x_1^2 + x_2^2)^2 & \text{if } \mathbf{x} \neq \mathbf{0} \\ 0 & \text{if } \mathbf{x} = \mathbf{0}. \end{cases} \end{aligned}$$

Note that

$$\frac{\partial f}{\partial x_1}([x_1, 0]^\top) = 0.$$

Hence,

$$\frac{\partial^2 f}{\partial x_1^2}(\mathbf{0}) = 0.$$

Also,

$$\frac{\partial f}{\partial x_1}([0, x_2]^\top) = -x_2.$$

Hence, the mixed partial is

$$\frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{0}) = -1.$$

We next compute

$$\frac{\partial^2 f}{\partial x_2^2} = \frac{\partial}{\partial x_2} \left( \frac{\partial f}{\partial x_2} \right),$$

where

$$\begin{aligned} \frac{\partial f}{\partial x_2}(x_1, x_2) \\ = \begin{cases} x_1(x_1^4 - x_2^4 - 4x_1^2x_2^2)/(x_1^2 + x_2^2)^2 & \text{if } \mathbf{x} \neq \mathbf{0} \\ 0 & \text{if } \mathbf{x} = \mathbf{0}. \end{cases} \end{aligned}$$

Note that

$$\frac{\partial f}{\partial x_2}([0, x_2]^\top) = 0.$$

Hence,

$$\frac{\partial^2 f}{\partial x_2^2}(\mathbf{0}) = 0.$$

Also,

$$\frac{\partial f}{\partial x_2}([x_1, 0]^\top) = x_1.$$

Hence, the mixed partial is

$$\frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{0}) = 1.$$

Therefore, the Hessian evaluated at the point  $\mathbf{0}$  is

$$\mathbf{F}(\mathbf{0}) = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix},$$

which is not symmetric. ■

## 5.4 Differentiation Rules

We now introduce the *chain rule* for differentiating the composition  $g(\mathbf{f}(t))$ , of a function  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^n$  and a function  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ .

**Theorem 5.6** *Let  $g : \mathcal{D} \rightarrow \mathbb{R}$  be differentiable on an open set  $\mathcal{D} \subset \mathbb{R}^n$ , and let  $\mathbf{f} : (a, b) \rightarrow \mathcal{D}$  be differentiable on  $(a, b)$ . Then, the composite function  $h : (a, b) \rightarrow \mathbb{R}$  given by  $h(t) = g(\mathbf{f}(t))$  is differentiable on  $(a, b)$ , and*

$$h'(t) = Dg(\mathbf{f}(t))D\mathbf{f}(t) = \nabla g(\mathbf{f}(t))^\top \begin{bmatrix} f'_1(t) \\ \vdots \\ f'_n(t) \end{bmatrix}.$$

□

*Proof.* By definition,

$$h'(t) = \lim_{s \rightarrow t} \frac{h(s) - h(t)}{s - t} = \lim_{s \rightarrow t} \frac{g(\mathbf{f}(s)) - g(\mathbf{f}(t))}{s - t}$$

if the limit exists. By Theorem 5.5 we write

$$g(\mathbf{f}(s)) - g(\mathbf{f}(t)) = Dg(\mathbf{f}(t))(\mathbf{f}(s) - \mathbf{f}(t)) + r(s),$$

where  $\lim_{s \rightarrow t} r(s)/(s - t) = 0$ . Therefore,

$$\frac{h(s) - h(t)}{s - t} = Dg(\mathbf{f}(t)) \frac{\mathbf{f}(s) - \mathbf{f}(t)}{s - t} + \frac{r(s)}{s - t}.$$

Letting  $s \rightarrow t$  yields

$$h'(t) = \lim_{s \rightarrow t} Dg(\mathbf{f}(t)) \frac{\mathbf{f}(s) - \mathbf{f}(t)}{s - t} + \frac{r(s)}{s - t} = Dg(\mathbf{f}(t))D\mathbf{f}(t).$$

■

Next, we present the *product rule*. Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be two differentiable functions. Define the function  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  by  $h(\mathbf{x}) = \mathbf{f}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$ . Then,  $h$  is also differentiable and

$$Dh(\mathbf{x}) = \mathbf{f}(\mathbf{x})^\top D\mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top D\mathbf{f}(\mathbf{x}).$$

We end this section with a list of some useful formulas from multivariable calculus. In each case, we compute the derivative with respect to  $\mathbf{x}$ . Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix and  $\mathbf{y} \in \mathbb{R}^m$  a given vector. Then,

$$\begin{aligned} D(\mathbf{y}^\top \mathbf{A}\mathbf{x}) &= \mathbf{y}^\top \mathbf{A} \\ D(\mathbf{x}^\top \mathbf{A}\mathbf{x}) &= \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top) \quad \text{if } m = n. \end{aligned}$$

It follows from the first formula above that if  $\mathbf{y} \in \mathbb{R}^n$ , then

$$D(\mathbf{y}^\top \mathbf{x}) = \mathbf{y}^\top.$$

It follows from the second formula above that if  $\mathbf{Q}$  is a symmetric matrix, then

$$D(\mathbf{x}^\top \mathbf{Q}\mathbf{x}) = 2\mathbf{x}^\top \mathbf{Q}.$$

In particular,

$$D(\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}^\top.$$

## 5.5 Level Sets and Gradients

The *level set* of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  at level  $c$  is the set of points

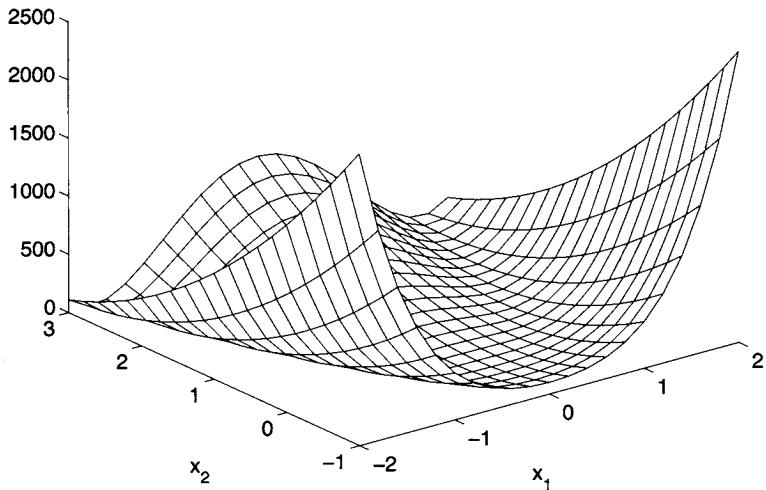
$$S = \{\mathbf{x} : f(\mathbf{x}) = c\}.$$

For  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , we are usually interested in  $S$  when it is a curve. For  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ , the sets  $S$  most often considered are surfaces.

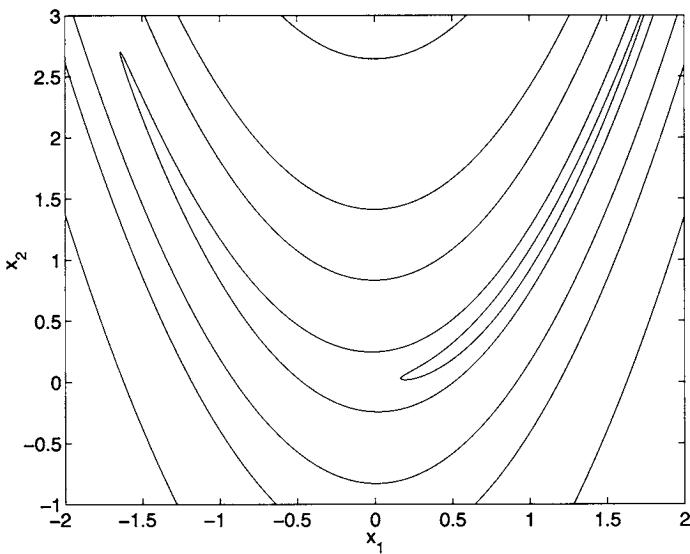
**Example 5.2** Consider the following real-valued function on  $\mathbb{R}^2$ :

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2, \quad \mathbf{x} = [x_1, x_2]^\top.$$

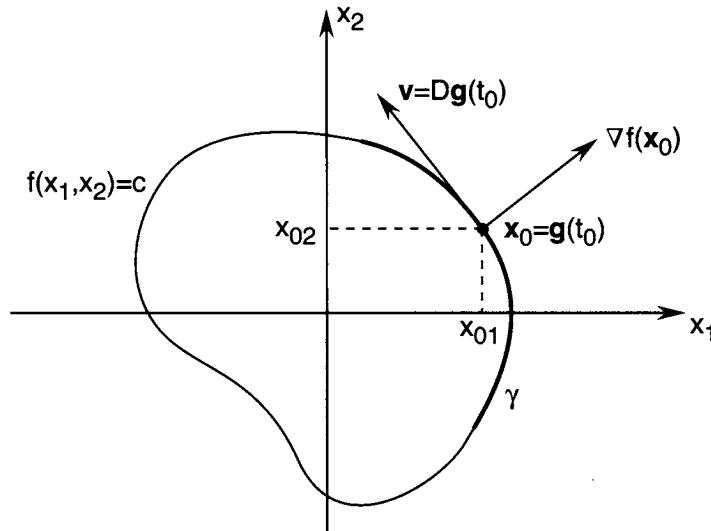
The function above is called *Rosenbrock's function*. A plot of the function  $f$  is shown in Figure 5.2. The level sets of  $f$  at levels 0.7, 7, 70, 200, and 700 are depicted in Figure 5.3. These level sets have a particular shape resembling



**Figure 5.2** Graph of Rosenbrock's function.



**Figure 5.3** Level sets of Rosenbrock's (banana) function.



**Figure 5.4** Orthogonality of the gradient to the level set.

bananas. For this reason, Rosenbrock's function is also called the *banana function*. ■

To say that a point  $\mathbf{x}_0$  is on the level set  $S$  at level  $c$  means that  $f(\mathbf{x}_0) = c$ . Now suppose that there is a curve  $\gamma$  lying in  $S$  and parameterized by a continuously differentiable function  $\mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^n$ . Suppose also that  $\mathbf{g}(t_0) = \mathbf{x}_0$  and  $D\mathbf{g}(t_0) = \mathbf{v} \neq \mathbf{0}$ , so that  $\mathbf{v}$  is a tangent vector to  $\gamma$  at  $\mathbf{x}_0$  (see Figure 5.4). Applying the chain rule to the function  $h(t) = f(\mathbf{g}(t))$  at  $t_0$  gives

$$h'(t_0) = Df(\mathbf{g}(t_0))D\mathbf{g}(t_0) = Df(\mathbf{x}_0)\mathbf{v}.$$

But since  $\gamma$  lies on  $S$ , we have

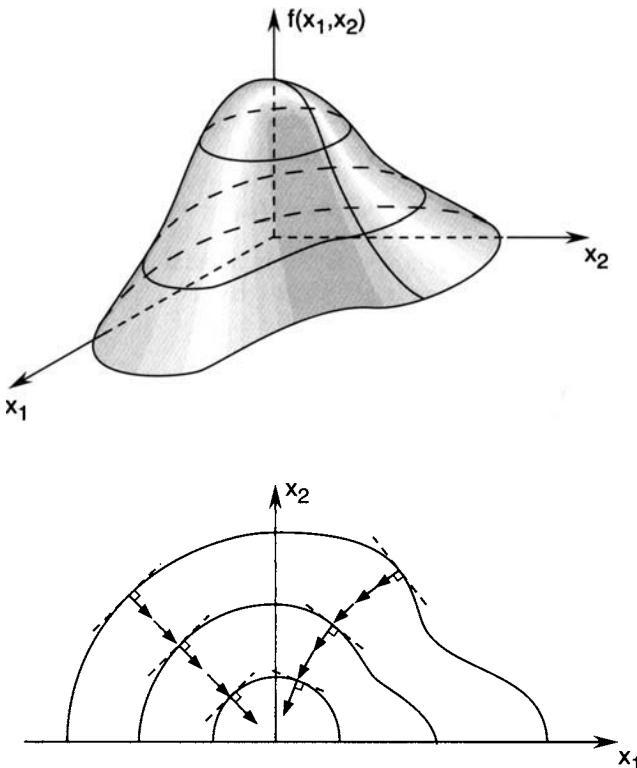
$$h(t) = f(\mathbf{g}(t)) = c;$$

that is,  $h$  is constant. Thus,  $h'(t_0) = 0$  and

$$Df(\mathbf{x}_0)\mathbf{v} = \nabla f(\mathbf{x}_0)^\top \mathbf{v} = 0.$$

Hence, we have proved, assuming  $f$  continuously differentiable, the following theorem (see Figure 5.4).

**Theorem 5.7** *The vector  $\nabla f(\mathbf{x}_0)$  is orthogonal to the tangent vector to an arbitrary smooth curve passing through  $\mathbf{x}_0$  on the level set determined by  $f(\mathbf{x}) = f(\mathbf{x}_0)$ .* □



**Figure 5.5** Illustration of a path of steepest ascent.

It is natural to say that  $\nabla f(\mathbf{x}_0)$  is *orthogonal* or *normal* to the level set  $S$  corresponding to  $\mathbf{x}_0$ , and it is also natural to take as the tangent plane (or line) to  $S$  at  $\mathbf{x}_0$  the set of all points  $\mathbf{x}$  satisfying

$$\nabla f(\mathbf{x}_0)^\top (\mathbf{x} - \mathbf{x}_0) = 0 \quad \text{if} \quad \nabla f(\mathbf{x}_0) \neq \mathbf{0}.$$

As we shall see later,  $\nabla f(\mathbf{x}_0)$  is the direction of *maximum rate of increase* of  $f$  at  $\mathbf{x}_0$ . Because  $\nabla f(\mathbf{x}_0)$  is orthogonal to the level set through  $\mathbf{x}_0$  determined by  $f(\mathbf{x}) = f(\mathbf{x}_0)$ , we deduce the following fact: The direction of maximum rate of increase of a real-valued differentiable function at a point is orthogonal to the level set of the function through that point.

Figure 5.5 illustrates the discussion above for the case  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ . The curve on the shaded surface in Figure 5.5 running from bottom to top has the property that its projection onto the  $(x_1, x_2)$ -plane is always orthogonal to the level curves and is called a *path of steepest ascent* because it always heads in the direction of maximum rate of increase for  $f$ .

The *graph* of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the set  $\{[\mathbf{x}^\top, f(\mathbf{x})]^\top : \mathbf{x} \in \mathbb{R}^n\} \subset \mathbb{R}^{n+1}$ . The notion of the gradient of a function has an alternative useful interpretation

in terms of the tangent hyperplane to its graph. To proceed, let  $\mathbf{x}_0 \in \mathbb{R}^n$  and  $z_0 = f(\mathbf{x}_0)$ . The point  $[\mathbf{x}_0^\top, z_0]^\top \in \mathbb{R}^{n+1}$  is a point on the graph of  $f$ . If  $f$  is differentiable at  $\xi$ , then the graph admits a nonvertical tangent hyperplane at  $\xi = [\mathbf{x}_0^\top, z_0]^\top$ . The hyperplane through  $\xi$  is the set of all points  $[x_1, \dots, x_n, z]^\top \in \mathbb{R}^{n+1}$  satisfying the equation

$$u_1(x_1 - x_{01}) + \dots + u_n(x_n - x_{0n}) + v(z - z_0) = 0,$$

where the vector  $[u_1, \dots, u_n, v]^\top \in \mathbb{R}^{n+1}$  is normal to the hyperplane. Assuming that this hyperplane is nonvertical (that is,  $v \neq 0$ ), let

$$d_i = -\frac{u_i}{v}.$$

Thus, we can rewrite the hyperplane equation above as

$$z = d_1(x_1 - x_{01}) + \dots + d_n(x_n - x_{0n}) + z_0.$$

We can think of the right side of the above equation as a function  $z : \mathbb{R}^n \rightarrow \mathbb{R}$ . Observe that for the hyperplane to be tangent to the graph of  $f$ , the functions  $f$  and  $z$  must have the same partial derivatives at the point  $\mathbf{x}_0$ . Hence, if  $f$  is differentiable at  $\mathbf{x}_0$ , its tangent hyperplane can be written in terms of its gradient, as given by the equation

$$z - z_0 = Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) = (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0).$$

## 5.6 Taylor Series

The basis for many numerical methods and models for optimization is Taylor's formula, which is given by Taylor's theorem.

**Theorem 5.8 Taylor's Theorem.** *Assume that a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is  $m$  times continuously differentiable (i.e.,  $f \in C^m$ ) on an interval  $[a, b]$ . Denote  $h = b - a$ . Then,*

$$f(b) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!} f^{(m-1)}(a) + R_m,$$

(called *Taylor's formula*) where  $f^{(i)}$  is the  $i$ th derivative of  $f$ , and

$$R_m = \frac{h^m (1-\theta)^{m-1}}{(m-1)!} f^{(m)}(a + \theta h) = \frac{h^m}{m!} f^{(m)}(a + \theta' h),$$

with  $\theta, \theta' \in (0, 1)$ . □

*Proof.* We have

$$R_m = f(b) - f(a) - \frac{h}{1!} f^{(1)}(a) - \frac{h^2}{2!} f^{(2)}(a) - \dots - \frac{h^{m-1}}{(m-1)!} f^{(m-1)}(a).$$

Denote by  $g_m(x)$  an auxiliary function obtained from  $R_m$  by replacing  $a$  by  $x$ . Hence,

$$\begin{aligned} g_m(x) &= f(b) - f(x) - \frac{b-x}{1!} f^{(1)}(x) - \frac{(b-x)^2}{2!} f^{(2)}(x) \\ &\quad - \dots - \frac{(b-x)^{m-1}}{(m-1)!} f^{(m-1)}(x). \end{aligned}$$

Differentiating  $g_m(x)$  yields

$$\begin{aligned} g_m^{(1)}(x) &= -f^{(1)}(x) + \left[ f^{(1)}(x) - \frac{b-x}{1!} f^{(2)}(x) \right] \\ &\quad + \left[ 2 \frac{b-x}{2!} f^{(2)}(x) - \frac{(b-x)^2}{2!} f^{(3)}(x) \right] + \dots \\ &\quad + \left[ (m-1) \frac{(b-x)^{m-2}}{(m-1)!} f^{(m-1)}(x) - \frac{(b-x)^{m-1}}{(m-1)!} f^{(m)}(x) \right] \\ &= -\frac{(b-x)^{m-1}}{(m-1)!} f^{(m)}(x). \end{aligned}$$

Observe that  $g_m(b) = 0$  and  $g_m(a) = R_m$ . Applying the mean-value theorem yields

$$\frac{g_m(b) - g_m(a)}{b-a} = g_m^{(1)}(a + \theta h),$$

where  $\theta \in (0, 1)$ . The equation above is equivalent to

$$-\frac{R_m}{h} = -\frac{(b-a-\theta h)^{m-1}}{(m-1)!} f^{(m)}(a+\theta h) = -\frac{h^{m-1}(1-\theta)^{m-1}}{(m-1)!} f^{(m)}(a+\theta h).$$

Hence,

$$R_m = \frac{h^m(1-\theta)^{m-1}}{(m-1)!} f^{(m)}(a+\theta h).$$

To derive the formula

$$R_m = \frac{h^m}{m!} f^{(m)}(a+\theta' h),$$

see, e.g., [81] or [83]. ■

An important property of Taylor's theorem arises from the form of the remainder  $R_m$ . To discuss this property further, we introduce the *order symbols*,  $O$  and  $o$ .

Let  $g$  be a real-valued function defined in some neighborhood of  $\mathbf{0} \in \mathbb{R}^n$ , with  $g(\mathbf{x}) \neq 0$  if  $\mathbf{x} \neq \mathbf{0}$ . Let  $\mathbf{f} : \Omega \rightarrow \mathbb{R}^m$  be defined in a domain  $\Omega \subset \mathbb{R}^n$  that includes  $\mathbf{0}$ . Then, we write

1.  $\mathbf{f}(\mathbf{x}) = O(g(\mathbf{x}))$  to mean that the quotient  $\|\mathbf{f}(\mathbf{x})\|/|g(\mathbf{x})|$  is bounded near 0; that is, there exist numbers  $K > 0$  and  $\delta > 0$  such that if  $\|\mathbf{x}\| < \delta$ ,  $\mathbf{x} \in \Omega$ , then  $\|\mathbf{f}(\mathbf{x})\|/|g(\mathbf{x})| \leq K$ .

2.  $\mathbf{f}(\mathbf{x}) = o(g(\mathbf{x}))$  to mean that

$$\lim_{\mathbf{x} \rightarrow 0, \mathbf{x} \in \Omega} \frac{\|\mathbf{f}(\mathbf{x})\|}{|g(\mathbf{x})|} = 0.$$

The symbol  $O(g(\mathbf{x}))$  [read “big-oh of  $g(\mathbf{x})$ ”] is used to represent a function that is bounded by a scaled version of  $g$  in a neighborhood of  $\mathbf{0}$ . Examples of such a function are:

- $x = O(x)$ .
- $\begin{bmatrix} x^3 \\ 2x^2 + 3x^4 \end{bmatrix} = O(x^2)$ .
- $\cos x = O(1)$ .
- $\sin x = O(x)$ .

On the other hand,  $o(g(\mathbf{x}))$  [read “little-oh of  $g(\mathbf{x})$ ”] represents a function that goes to zero “faster” than  $g(\mathbf{x})$  in the sense that  $\lim_{\mathbf{x} \rightarrow \mathbf{0}} \|o(g(\mathbf{x}))\|/|g(\mathbf{x})| = 0$ . Examples of such functions are:

- $x^2 = o(x)$ .
- $\begin{bmatrix} x^3 \\ 2x^2 + 3x^4 \end{bmatrix} = o(x)$ .
- $x^3 = o(x^2)$ .
- $x = o(1)$ .

Note that if  $\mathbf{f}(\mathbf{x}) = o(g(\mathbf{x}))$ , then  $\mathbf{f}(\mathbf{x}) = O(g(\mathbf{x}))$  (but the converse is not necessarily true). Also, if  $\mathbf{f}(\mathbf{x}) = O(\|\mathbf{x}\|^p)$ , then  $\mathbf{f}(\mathbf{x}) = o(\|\mathbf{x}\|^{p-\varepsilon})$  for any  $\varepsilon > 0$ .

Suppose that  $f \in \mathcal{C}^m$ . Recall that the remainder term in Taylor’s theorem has the form

$$R_m = \frac{h^m}{m!} f^{(m)}(a + \theta h),$$

where  $\theta \in (0, 1)$ . Substituting this into Taylor’s formula, we get

$$f(b) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^{m-1}}{(m-1)!} f^{(m-1)}(a) + \frac{h^m}{m!} f^{(m)}(a + \theta h).$$

By the continuity of  $f^{(m)}$ , we have  $f^{(m)}(a + \theta h) \rightarrow f^{(m)}(a)$  as  $h \rightarrow 0$ ; that is,  $f^{(m)}(a + \theta h) = f^{(m)}(a) + o(1)$ . Therefore,

$$\frac{h^m}{m!} f^{(m)}(a + \theta h) = \frac{h^m}{m!} f^{(m)}(a) + o(h^m),$$

since  $h^m o(1) = o(h^m)$ . We may then write Taylor's formula as

$$f(b) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + o(h^m).$$

If, in addition, we assume that  $f \in \mathcal{C}^{m+1}$ , we may replace the term  $o(h^m)$  above by  $O(h^{m+1})$ . To see this, we first write Taylor's formula with  $R_{m+1}$ :

$$f(b) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + R_{m+1},$$

where

$$R_{m+1} = \frac{h^{m+1}}{(m+1)!} f^{(m+1)}(a + \theta' h),$$

with  $\theta' \in (0, 1)$ . Because  $f^{(m+1)}$  is bounded on  $[a, b]$  (by Theorem 4.2),

$$R_{m+1} = O(h^{m+1}).$$

Therefore, if  $f \in \mathcal{C}^{m+1}$ , we may write Taylor's formula as

$$f(b) = f(a) + \frac{h}{1!} f^{(1)}(a) + \frac{h^2}{2!} f^{(2)}(a) + \cdots + \frac{h^m}{m!} f^{(m)}(a) + O(h^{m+1}).$$

We now turn to the Taylor series expansion of a real-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  about the point  $\mathbf{x}_0 \in \mathbb{R}^n$ . Suppose that  $f \in \mathcal{C}^2$ . Let  $\mathbf{x}$  and  $\mathbf{x}_0$  be points in  $\mathbb{R}^n$ , and let  $\mathbf{z}(\alpha) = \mathbf{x}_0 + \alpha(\mathbf{x} - \mathbf{x}_0)/\|\mathbf{x} - \mathbf{x}_0\|$ . Define  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  by

$$\phi(\alpha) = f(\mathbf{z}(\alpha)) = f(\mathbf{x}_0 + \alpha(\mathbf{x} - \mathbf{x}_0)/\|\mathbf{x} - \mathbf{x}_0\|).$$

Using the chain rule, we obtain

$$\begin{aligned} \phi'(\alpha) &= \frac{d\phi}{d\alpha}(\alpha) \\ &= Df(\mathbf{z}(\alpha))D\mathbf{z}(\alpha) = Df(\mathbf{z}(\alpha))\frac{(\mathbf{x} - \mathbf{x}_0)}{\|\mathbf{x} - \mathbf{x}_0\|} \\ &= (\mathbf{x} - \mathbf{x}_0)^\top Df(\mathbf{z}(\alpha))^\top / \|\mathbf{x} - \mathbf{x}_0\| \end{aligned}$$

and

$$\begin{aligned}
\phi''(\alpha) &= \frac{d^2\phi}{d\alpha^2}(\alpha) \\
&= \frac{d}{d\alpha} \left( \frac{d\phi}{d\alpha} \right)(\alpha) \\
&= \frac{(\mathbf{x} - \mathbf{x}_0)^\top}{\|\mathbf{x} - \mathbf{x}_0\|} \frac{d}{d\alpha} Df(\mathbf{z}(\alpha))^\top \\
&= \frac{(\mathbf{x} - \mathbf{x}_0)^\top}{\|\mathbf{x} - \mathbf{x}_0\|} D(Df)(\mathbf{z}(\alpha))^\top \frac{d\mathbf{z}}{d\alpha}(\alpha) \\
&= \frac{1}{\|\mathbf{x} - \mathbf{x}_0\|^2} (\mathbf{x} - \mathbf{x}_0)^\top D^2 f(\mathbf{z}(\alpha))^\top (\mathbf{x} - \mathbf{x}_0) \\
&= \frac{1}{\|\mathbf{x} - \mathbf{x}_0\|^2} (\mathbf{x} - \mathbf{x}_0)^\top D^2 f(\mathbf{z}(\alpha))(\mathbf{x} - \mathbf{x}_0),
\end{aligned}$$

where we recall that

$$D^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}.$$

Observe that

$$\begin{aligned}
f(\mathbf{x}) &= \phi(\|\mathbf{x} - \mathbf{x}_0\|) \\
&= \phi(0) + \frac{\|\mathbf{x} - \mathbf{x}_0\|}{1!} \phi'(0) + \frac{\|\mathbf{x} - \mathbf{x}_0\|^2}{2!} \phi''(0) + o(\|\mathbf{x} - \mathbf{x}_0\|^2).
\end{aligned}$$

Hence,

$$\begin{aligned}
f(\mathbf{x}) &= f(\mathbf{x}_0) + \frac{1}{1!} Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \\
&\quad + \frac{1}{2!} (\mathbf{x} - \mathbf{x}_0)^\top D^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|^2).
\end{aligned}$$

If we assume that  $f \in \mathcal{C}^3$ , we may use the formula for the remainder term  $R_3$  to conclude that

$$\begin{aligned}
f(\mathbf{x}) &= f(\mathbf{x}_0) + \frac{1}{1!} Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \\
&\quad + \frac{1}{2!} (\mathbf{x} - \mathbf{x}_0)^\top D^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + O(\|\mathbf{x} - \mathbf{x}_0\|^3).
\end{aligned}$$

We end with a statement of the *mean value theorem*, which is closely related to Taylor's theorem.

**Theorem 5.9** If a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable on an open set  $\Omega \subset \mathbb{R}^n$ , then for any pair of points  $\mathbf{x}, \mathbf{y} \in \Omega$ , there exists a matrix  $\mathbf{M}$  such that

$$\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y}) = \mathbf{M}(\mathbf{x} - \mathbf{y}).$$

□

The mean value theorem follows from Taylor's theorem (for the case where  $m = 1$ ) applied to each component of  $\mathbf{f}$ . It is easy to see that  $\mathbf{M}$  is a matrix whose rows are the rows of  $D\mathbf{f}$  evaluated at points that lie on the line segment joining  $\mathbf{x}$  and  $\mathbf{y}$  (these points may differ from row to row).

For further reading in calculus, consult [13], [81], [83], [115], [120], [134]. A basic treatment of real analysis can be found in [2], [112], whereas a more advanced treatment is provided in [89], [111]. For stimulating reading on the “big-oh” notation, see [77, pp. 104–108].

## EXERCISES

**5.1** Show that a sufficient condition for  $\lim_{k \rightarrow \infty} \mathbf{A}^k = \mathbf{O}$  is  $\|\mathbf{A}\| < 1$ .

**5.2** Show that for any matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,

$$\|\mathbf{A}\| \geq \max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})|.$$

*Hint:* Use Exercise 5.1.

**5.3** Consider the function

$$f(\mathbf{x}) = (\mathbf{a}^\top \mathbf{x})(\mathbf{b}^\top \mathbf{x}),$$

where  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{x}$  are  $n$ -dimensional vectors.

- a. Find  $\nabla f(\mathbf{x})$ .
- b. Find the Hessian  $\mathbf{F}(\mathbf{x})$ .

**5.4** Define the functions  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $g : \mathbb{R} \rightarrow \mathbb{R}^2$  by  $f(\mathbf{x}) = x_1^2/6 + x_2^2/4$ ,  $\mathbf{g}(t) = [3t + 5, 2t - 6]^\top$ . Let  $F : \mathbb{R} \rightarrow \mathbb{R}$  be given by  $F(t) = f(g(t))$ . Evaluate  $\frac{dF}{dt}(t)$  using the chain rule.

**5.5** Consider  $f(\mathbf{x}) = x_1x_2/2$ ,  $\mathbf{g}(s, t) = [4s + 3t, 2s + t]^\top$ . Evaluate  $\frac{\partial}{\partial s} f(\mathbf{g}(s, t))$  and  $\frac{\partial}{\partial t} f(\mathbf{g}(s, t))$  using the chain rule.

**5.6** Let  $\mathbf{x}(t) = [e^t + t^3, t^2, t + 1]^\top$ ,  $t \in \mathbb{R}$ , and  $f(\mathbf{x}) = x_1^3x_2x_3^2 + x_1x_2 + x_3$ ,  $\mathbf{x} = [x_1, x_2, x_3]^\top \in \mathbb{R}^3$ . Find  $\frac{d}{dt} f(\mathbf{x}(t))$  in terms of  $t$ .

**5.7** Suppose that  $\mathbf{f}(\mathbf{x}) = o(g(\mathbf{x}))$ . Show that for any given  $\varepsilon > 0$ , there exists  $\delta > 0$  such that if  $\|\mathbf{x}\| < \delta$ , then  $\|\mathbf{f}(\mathbf{x})\| < \varepsilon|g(\mathbf{x})|$ .

**5.8** Use Exercise 5.7 to show that if functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  satisfy  $f(\mathbf{x}) = -g(\mathbf{x}) + o(g(\mathbf{x}))$  and  $g(\mathbf{x}) > 0$  for all  $\mathbf{x} \neq \mathbf{0}$ , then for all  $\mathbf{x} \neq \mathbf{0}$  sufficiently small, we have  $f(\mathbf{x}) < 0$ .

**5.9** Let

$$\begin{aligned} f_1(x_1, x_2) &= x_1^2 - x_2^2, \\ f_2(x_1, x_2) &= 2x_1x_2. \end{aligned}$$

Sketch the level sets associated with  $f_1(x_1, x_2) = 12$  and  $f_2(x_1, x_2) = 16$  on the same diagram. Indicate on the diagram the values of  $\mathbf{x} = [x_1, x_2]^\top$  for which  $\mathbf{f}(\mathbf{x}) = [f_1(x_1, x_2), f_2(x_1, x_2)]^\top = [12, 16]^\top$ .

**5.10** Write down the Taylor series expansion of the following functions about the given points  $\mathbf{x}_0$ . Neglect terms of order three or higher.

- a.  $f(\mathbf{x}) = x_1 e^{-x_2} + x_2 + 1$ ,  $\mathbf{x}_0 = [1, 0]^\top$ .
- b.  $f(\mathbf{x}) = x_1^4 + 2x_1^2x_2^2 + x_2^4$ ,  $\mathbf{x}_0 = [1, 1]^\top$ .
- c.  $f(\mathbf{x}) = e^{x_1-x_2} + e^{x_1+x_2} + x_1 + x_2 + 1$ ,  $\mathbf{x}_0 = [1, 0]^\top$ .

## **PART II**

---

# **UNCONSTRAINED OPTIMIZATION**

---



# CHAPTER 6

---

## BASICS OF SET-CONSTRAINED AND UNCONSTRAINED OPTIMIZATION

---

### 6.1 Introduction

In this chapter we consider the optimization problem

$$\begin{aligned} & \text{minimize} && f(\boldsymbol{x}) \\ & \text{subject to} && \boldsymbol{x} \in \Omega. \end{aligned}$$

The function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that we wish to minimize is a real-valued function called the *objective function* or *cost function*. The vector  $\boldsymbol{x}$  is an  $n$ -vector of independent variables:  $\boldsymbol{x} = [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^n$ . The variables  $x_1, \dots, x_n$  are often referred to as *decision variables*. The set  $\Omega$  is a subset of  $\mathbb{R}^n$  called the *constraint set* or *feasible set*.

The optimization problem above can be viewed as a decision problem that involves finding the “best” vector  $\boldsymbol{x}$  of the decision variables over all possible vectors in  $\Omega$ . By the “best” vector we mean the one that results in the smallest value of the objective function. This vector is called the *minimizer* of  $f$  over  $\Omega$ . It is possible that there may be many minimizers. In this case, finding any of the minimizers will suffice.

There are also optimization problems that require maximization of the objective function, in which case we seek *maximizers*. Minimizers and maximizers are also called *extremizers*. Maximization problems, however, can be represented equivalently in the minimization form above because maximizing  $f$  is equivalent to minimizing  $-f$ . Therefore, we can confine our attention to minimization problems without loss of generality.

The problem above is a general form of a *constrained optimization problem*, because the decision variables are constrained to be in the constraint set  $\Omega$ . If  $\Omega = \mathbb{R}^n$ , then we refer to the problem as an *unconstrained optimization problem*. In this chapter we discuss basic properties of the general optimization problem above, which includes the unconstrained case. In the remaining chapters of this part, we deal with iterative algorithms for solving unconstrained optimization problems.

The constraint " $\mathbf{x} \in \Omega$ " is called a *set constraint*. Often, the constraint set  $\Omega$  takes the form  $\Omega = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ , where  $\mathbf{h}$  and  $\mathbf{g}$  are given functions. We refer to such constraints as *functional constraints*. The remainder of this chapter deals with general set constraints, including the special case where  $\Omega = \mathbb{R}^n$ . The case where  $\Omega = \mathbb{R}^n$  is called the *unconstrained* case. In Parts III and IV we consider constrained optimization problems with functional constraints.

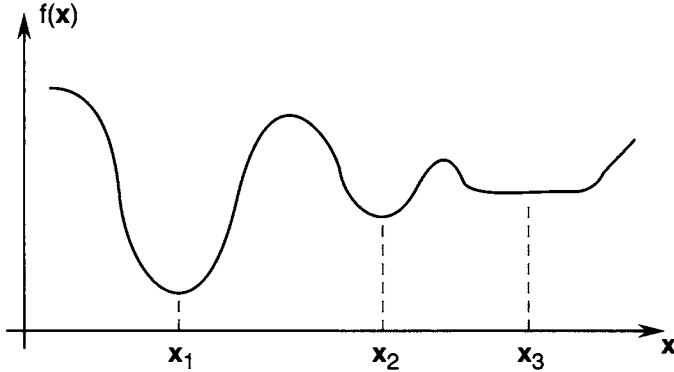
In considering the general optimization problem above, we distinguish between two kinds of minimizers, as specified by the following definitions.

**Definition 6.1** Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a real-valued function defined on some set  $\Omega \subset \mathbb{R}^n$ . A point  $\mathbf{x}^* \in \Omega$  is a *local minimizer* of  $f$  over  $\Omega$  if there exists  $\varepsilon > 0$  such that  $f(\mathbf{x}) \geq f(\mathbf{x}^*)$  for all  $\mathbf{x} \in \Omega \setminus \{\mathbf{x}^*\}$  and  $\|\mathbf{x} - \mathbf{x}^*\| < \varepsilon$ . A point  $\mathbf{x}^* \in \Omega$  is a *global minimizer* of  $f$  over  $\Omega$  if  $f(\mathbf{x}) \geq f(\mathbf{x}^*)$  for all  $\mathbf{x} \in \Omega \setminus \{\mathbf{x}^*\}$ . ■

If in the definitions above we replace " $\geq$ " with " $>$ ," then we have a *strict local minimizer* and a *strict global minimizer*, respectively. In Figure 6.1, we illustrate the definitions for  $n = 1$ .

If  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ , we write  $f(\mathbf{x}^*) = \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$  and  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \Omega} f(\mathbf{x})$ . If the minimization is unconstrained, we simply write  $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$  or  $\mathbf{x}^* = \arg \min f(\mathbf{x})$ . In other words, given a real-valued function  $f$ , the notation  $\arg \min f(\mathbf{x})$  denotes the *argument* that minimizes the function  $f$  (a point in the domain of  $f$ ), assuming that such a point is unique (if there is more than one such point, we pick one arbitrarily). For example, if  $f : \mathbb{R} \rightarrow \mathbb{R}$  is given by  $f(x) = (x+1)^2 + 3$ , then  $\arg \min f(x) = -1$ . If we write  $\arg \min_{x \in \Omega}$ , then we treat " $x \in \Omega$ " to be a constraint for the minimization. For example, for the function  $f$  above,  $\arg \min_{x \geq 0} f(x) = 0$ .

Strictly speaking, an optimization problem is solved only when a global minimizer is found. However, global minimizers are, in general, difficult to find. Therefore, in practice, we often have to be satisfied with finding local minimizers.



**Figure 6.1** Examples of minimizers:  $\mathbf{x}_1$ : strict global minimizer;  $\mathbf{x}_2$ : strict local minimizer;  $\mathbf{x}_3$ : local (not strict) minimizer.

## 6.2 Conditions for Local Minimizers

In this section we derive conditions for a point  $\mathbf{x}^*$  to be a local minimizer. We use derivatives of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Recall that the first-order derivative of  $f$ , denoted  $Df$ , is

$$Df \triangleq \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right].$$

Note that the gradient  $\nabla f$  is just the transpose of  $Df$ ; that is,  $\nabla f = (Df)^\top$ . The second derivative of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  (also called the *Hessian* of  $f$ ) is

$$\mathbf{F}(\mathbf{x}) \triangleq D^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}.$$

**Example 6.1** Let  $f(x_1, x_2) = 5x_1 + 8x_2 + x_1x_2 - x_1^2 - 2x_2^2$ . Then,

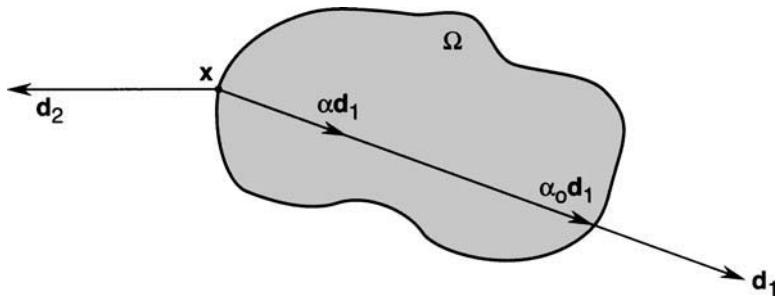
$$Df(\mathbf{x}) = (\nabla f(\mathbf{x}))^\top = \left[ \frac{\partial f}{\partial x_1}(\mathbf{x}), \frac{\partial f}{\partial x_2}(\mathbf{x}) \right] = [5 + x_2 - 2x_1, 8 + x_1 - 4x_2]$$

and

$$\mathbf{F}(\mathbf{x}) = D^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} -2 & 1 \\ 1 & -4 \end{bmatrix}.$$

■

Given an optimization problem with constraint set  $\Omega$ , a minimizer may lie either in the interior or on the boundary of  $\Omega$ . To study the case where it lies on the boundary, we need the notion of *feasible directions*.



**Figure 6.2** Two-dimensional illustration of feasible directions;  $\mathbf{d}_1$  is a feasible direction,  $\mathbf{d}_2$  is not a feasible direction.

**Definition 6.2** A vector  $\mathbf{d} \in \mathbb{R}^n$ ,  $\mathbf{d} \neq \mathbf{0}$ , is a *feasible direction* at  $\mathbf{x} \in \Omega$  if there exists  $\alpha_0 > 0$  such that  $\mathbf{x} + \alpha\mathbf{d} \in \Omega$  for all  $\alpha \in [0, \alpha_0]$ . ■

Figure 6.2 illustrates the notion of feasible directions.

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a real-valued function and let  $\mathbf{d}$  be a feasible direction at  $\mathbf{x} \in \Omega$ . The *directional derivative of  $f$  in the direction  $\mathbf{d}$* , denoted  $\partial f / \partial \mathbf{d}$ , is the real-valued function defined by

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}) = \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha\mathbf{d}) - f(\mathbf{x})}{\alpha}.$$

If  $\|\mathbf{d}\| = 1$ , then  $\partial f / \partial \mathbf{d}$  is the rate of increase of  $f$  at  $\mathbf{x}$  in the direction  $\mathbf{d}$ . To compute the directional derivative above, suppose that  $\mathbf{x}$  and  $\mathbf{d}$  are given. Then,  $f(\mathbf{x} + \alpha\mathbf{d})$  is a function of  $\alpha$ , and

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}) = \frac{d}{d\alpha} f(\mathbf{x} + \alpha\mathbf{d}) \Big|_{\alpha=0}.$$

Applying the chain rule yields

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}) = \frac{d}{d\alpha} f(\mathbf{x} + \alpha\mathbf{d}) \Big|_{\alpha=0} = \nabla f(\mathbf{x})^\top \mathbf{d} = \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle = \mathbf{d}^\top \nabla f(\mathbf{x}).$$

In summary, if  $\mathbf{d}$  is a unit vector ( $\|\mathbf{d}\| = 1$ ), then  $\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle$  is the rate of increase of  $f$  at the point  $\mathbf{x}$  in the direction  $\mathbf{d}$ .

**Example 6.2** Define  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  by  $f(\mathbf{x}) = x_1x_2x_3$ , and let

$$\mathbf{d} = \left[ \frac{1}{2}, \frac{1}{2}, \frac{1}{\sqrt{2}} \right]^\top.$$

The directional derivative of  $f$  in the direction  $\mathbf{d}$  is

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}) = \nabla f(\mathbf{x})^\top \mathbf{d} = [x_2x_3, x_1x_3, x_1x_2] \begin{bmatrix} 1/2 \\ 1/2 \\ 1/\sqrt{2} \end{bmatrix} = \frac{x_2x_3 + x_1x_3 + \sqrt{2}x_1x_2}{2}.$$

Note that because  $\|\mathbf{d}\| = 1$ , the above is also the rate of increase of  $f$  at  $\mathbf{x}$  in the direction  $\mathbf{d}$ . ■

We are now ready to state and prove the following theorem.

**Theorem 6.1 First-Order Necessary Condition (FONC).** *Let  $\Omega$  be a subset of  $\mathbb{R}^n$  and  $f \in C^1$  a real-valued function on  $\Omega$ . If  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega$ , then for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have*

$$\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0.$$

□

*Proof.* Define

$$\mathbf{x}(\alpha) = \mathbf{x}^* + \alpha\mathbf{d} \in \Omega.$$

Note that  $\mathbf{x}(0) = \mathbf{x}^*$ . Define the composite function

$$\phi(\alpha) = f(\mathbf{x}(\alpha)).$$

Then, by Taylor's theorem,

$$f(\mathbf{x}^* + \alpha\mathbf{d}) - f(\mathbf{x}^*) = \phi(\alpha) - \phi(0) = \phi'(0)\alpha + o(\alpha) = \alpha\mathbf{d}^\top \nabla f(\mathbf{x}(0)) + o(\alpha),$$

where  $\alpha \geq 0$  [recall the definition of  $o(\alpha)$  ("little-oh of  $\alpha$ ") in Part I]. Thus, if  $\phi(\alpha) \geq \phi(0)$ , that is,  $f(\mathbf{x}^* + \alpha\mathbf{d}) \geq f(\mathbf{x}^*)$  for sufficiently small values of  $\alpha > 0$  ( $\mathbf{x}^*$  is a local minimizer), then we have to have  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0$  (see Exercise 5.8). ■

Theorem 6.1 is illustrated in Figure 6.3.

An alternative way to express the FONC is

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}^*) \geq 0$$

for all feasible directions  $\mathbf{d}$ . In other words, if  $\mathbf{x}^*$  is a local minimizer, then the rate of increase of  $f$  at  $\mathbf{x}^*$  in any feasible direction  $\mathbf{d}$  in  $\Omega$  is nonnegative. Using directional derivatives, an alternative proof of Theorem 6.1 is as follows. Suppose that  $\mathbf{x}^*$  is a local minimizer. Then, for any feasible direction  $\mathbf{d}$ , there exists  $\bar{\alpha} > 0$  such that for all  $\alpha \in (0, \bar{\alpha})$ ,

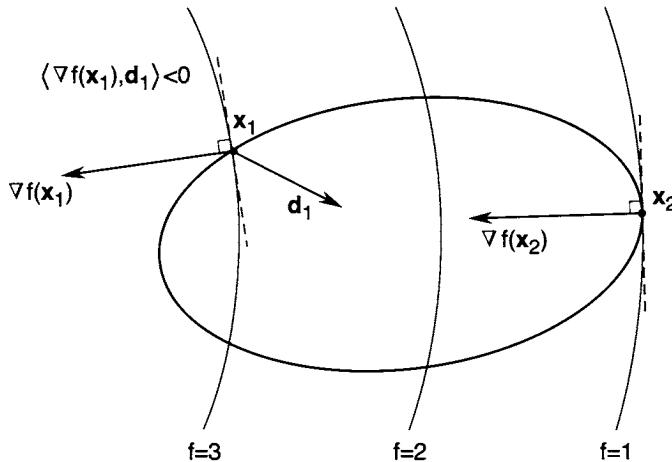
$$f(\mathbf{x}^*) \leq f(\mathbf{x}^* + \alpha\mathbf{d}).$$

Hence, for all  $\alpha \in (0, \bar{\alpha})$ , we have

$$\frac{f(\mathbf{x}^* + \alpha\mathbf{d}) - f(\mathbf{x}^*)}{\alpha} \geq 0.$$

Taking the limit as  $\alpha \rightarrow 0$ , we conclude that

$$\frac{\partial f}{\partial \mathbf{d}}(\mathbf{x}^*) \geq 0.$$



**Figure 6.3** Illustration of the FONC for a constrained case;  $x_1$  does not satisfy the FONC, whereas  $x_2$  satisfies the FONC.

A special case of interest is when  $\mathbf{x}^*$  is an interior point of  $\Omega$  (see Section 4.4). In this case, any direction is feasible, and we have the following result.

**Corollary 6.1 Interior Case.** Let  $\Omega$  be a subset of  $\mathbb{R}^n$  and  $f \in \mathcal{C}^1$  a real-valued function on  $\Omega$ . If  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega$  and if  $\mathbf{x}^*$  is an interior point of  $\Omega$ , then

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

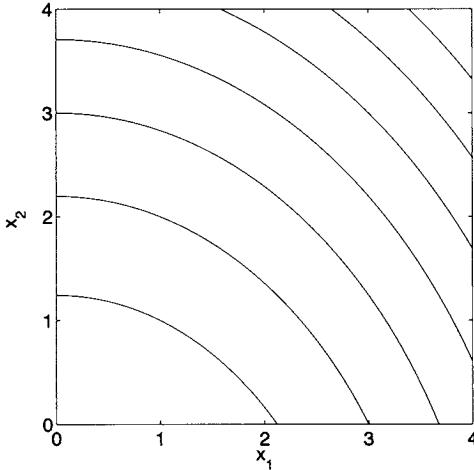
□

*Proof.* Suppose that  $f$  has a local minimizer  $\mathbf{x}^*$  that is an interior point of  $\Omega$ . Because  $\mathbf{x}^*$  is an interior point of  $\Omega$ , the set of feasible directions at  $\mathbf{x}^*$  is the whole of  $\mathbb{R}^n$ . Thus, for any  $\mathbf{d} \in \mathbb{R}^n$ ,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0$  and  $-\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0$ . Hence,  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$  for all  $\mathbf{d} \in \mathbb{R}^n$ , which implies that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . ■

**Example 6.3** Consider the problem

$$\begin{aligned} \text{minimize } & x_1^2 + 0.5x_2^2 + 3x_2 + 4.5 \\ \text{subject to } & x_1, x_2 \geq 0. \end{aligned}$$

- a. Is the first-order necessary condition (FONC) for a local minimizer satisfied at  $\mathbf{x} = [1, 3]^\top$ ?
- b. Is the FONC for a local minimizer satisfied at  $\mathbf{x} = [0, 3]^\top$ ?
- c. Is the FONC for a local minimizer satisfied at  $\mathbf{x} = [1, 0]^\top$ ?



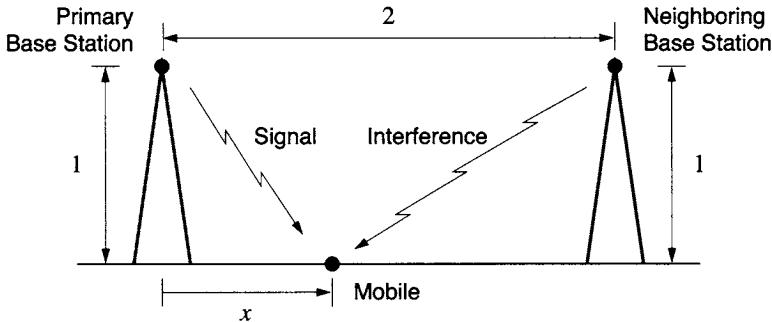
**Figure 6.4** Level sets of the function in Example 6.3.

- d. Is the FONC for a local minimizer satisfied at  $\mathbf{x} = [0, 0]^\top$ ?

*Solution:* First, let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be defined by  $f(\mathbf{x}) = x_1^2 + 0.5x_2^2 + 3x_2 + 4.5$ , where  $\mathbf{x} = [x_1, x_2]^\top$ . A plot of the level sets of  $f$  is shown in Figure 6.4.

- a. At  $\mathbf{x} = [1, 3]^\top$ , we have  $\nabla f(\mathbf{x}) = [2x_1, x_2 + 3]^\top = [2, 6]^\top$ . The point  $\mathbf{x} = [1, 3]^\top$  is an interior point of  $\Omega = \{\mathbf{x} : x_1 \geq 0, x_2 \geq 0\}$ . Hence, the FONC requires that  $\nabla f(\mathbf{x}) = \mathbf{0}$ . The point  $\mathbf{x} = [1, 3]^\top$  does not satisfy the FONC for a local minimizer.
- b. At  $\mathbf{x} = [0, 3]^\top$ , we have  $\nabla f(\mathbf{x}) = [0, 6]^\top$ , and hence  $\mathbf{d}^\top \nabla f(\mathbf{x}) = 6d_2$ , where  $\mathbf{d} = [d_1, d_2]^\top$ . For  $\mathbf{d}$  to be feasible at  $\mathbf{x}$ , we need  $d_1 \geq 0$ , and  $d_2$  can take an arbitrary value in  $\mathbb{R}$ . The point  $\mathbf{x} = [0, 3]^\top$  does not satisfy the FONC for a minimizer because  $d_2$  is allowed to be less than zero. For example,  $\mathbf{d} = [1, -1]^\top$  is a feasible direction, but  $\mathbf{d}^\top \nabla f(\mathbf{x}) = -6 < 0$ .
- c. At  $\mathbf{x} = [1, 0]^\top$ , we have  $\nabla f(\mathbf{x}) = [2, 3]^\top$ , and hence  $\mathbf{d}^\top \nabla f(\mathbf{x}) = 2d_1 + 3d_2$ . For  $\mathbf{d}$  to be feasible, we need  $d_2 \geq 0$ , and  $d_1$  can take an arbitrary value in  $\mathbb{R}$ . For example,  $\mathbf{d} = [-5, 1]^\top$  is a feasible direction. But  $\mathbf{d}^\top \nabla f(\mathbf{x}) = -7 < 0$ . Thus,  $\mathbf{x} = [1, 0]^\top$  does not satisfy the FONC for a local minimizer.
- d. At  $\mathbf{x} = [0, 0]^\top$ , we have  $\nabla f(\mathbf{x}) = [0, 3]^\top$ , and hence  $\mathbf{d}^\top \nabla f(\mathbf{x}) = 3d_2$ . For  $\mathbf{d}$  to be feasible, we need  $d_2 \geq 0$  and  $d_1 \geq 0$ . Hence,  $\mathbf{x} = [0, 0]^\top$  satisfies the FONC for a local minimizer. ■

**Example 6.4** Figure 6.5 shows a simplified model of a cellular wireless system (the distances shown have been scaled down to make the calculations



**Figure 6.5** Simplified cellular wireless system in Example 6.4.

simpler). A mobile user (also called a *mobile*) is located at position  $x$  (see Figure 6.5).

There are two base station antennas, one for the primary base station and another for the neighboring base station. Both antennas are transmitting signals to the mobile user, at equal power. However, the power of the received signal as measured by the mobile is the reciprocal of the squared distance from the associated antenna (primary or neighboring base station). We are interested in finding the position of the mobile that maximizes the *signal-to-interference ratio*, which is the ratio of the signal power received from the primary base station to the signal power received from the neighboring base station.

We use the FONC to solve this problem. The squared distance from the mobile to the primary antenna is  $1 + x^2$ , while the squared distance from the mobile to the neighboring antenna is  $1 + (2 - x)^2$ . Therefore, the signal-to-interference ratio is

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

We have

$$\begin{aligned} f'(x) &= \frac{-2(2-x)(1+x^2) - 2x(1+(2-x)^2)}{(1+x^2)^2} \\ &= \frac{4(x^2 - 2x - 1)}{(1+x^2)^2}. \end{aligned}$$

By the FONC, at the optimal position  $x^*$  we have  $f'(x^*) = 0$ . Hence, either  $x^* = 1 - \sqrt{2}$  or  $x^* = 1 + \sqrt{2}$ . Evaluating the objective function at these two candidate points, it is easy to see that  $x^* = 1 - \sqrt{2}$  is the optimal position. ■

The next example illustrates that in some problems the FONC is not helpful for eliminating candidate local minimizers. However, in such cases, there may be a recasting of the problem into an equivalent form that makes the FONC useful.

**Example 6.5** Consider the set-constrained problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $\Omega = \{[\mathbf{x}_1, \mathbf{x}_2]^\top : \mathbf{x}_1^2 + \mathbf{x}_2^2 = 1\}$ .

- a. Consider a point  $\mathbf{x}^* \in \Omega$ . Specify all feasible directions at  $\mathbf{x}^*$ .
- b. Which points in  $\Omega$  satisfy the FONC for this set-constrained problem?
- c. Based on part b, is the FONC for this set-constrained problem useful for eliminating local-minimizer candidates?
- d. Suppose that we use polar coordinates to parameterize points  $\mathbf{x} \in \Omega$  in terms of a single parameter  $\theta$ :

$$x_1 = \cos \theta \quad x_2 = \sin \theta.$$

Now use the FONC for unconstrained problems (with respect to  $\theta$ ) to derive a necessary condition of this sort: If  $\mathbf{x}^* \in \Omega$  is a local minimizer, then  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$  for all  $\mathbf{d}$  satisfying a “certain condition.” Specify what this certain condition is.

*Solution:*

- a. There are no feasible directions at any  $\mathbf{x}^*$ .
- b. Because of part a, *all* points in  $\Omega$  satisfy the FONC for this set-constrained problem.
- c. No, the FONC for this set-constrained problem is not useful for eliminating local-minimizer candidates.
- d. Write  $h(\theta) = f(g(\theta))$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}^2$  is given by the equations relating  $\theta$  to  $\mathbf{x} = [x_1, x_2]^\top$ . Note that  $Dg(\theta) = [-\sin \theta, \cos \theta]^\top$ . Hence, by the chain rule,

$$h'(\theta) = Df(g(\theta))Dg(\theta) = Dg(\theta)^\top \nabla f(g(\theta)).$$

Notice that  $Dg(\theta)$  is tangent to  $\Omega$  at  $\mathbf{x} = g(\theta)$ . Alternatively, we could say that  $Dg(\theta)$  is orthogonal to  $\mathbf{x} = g(\theta)$ .

Suppose that  $\mathbf{x}^* \in \Omega$  is a local minimizer. Write  $\mathbf{x}^* = g(\theta^*)$ . Then  $\theta^*$  is an unconstrained minimizer of  $h$ . By the FONC for unconstrained problems,  $h'(\theta^*) = 0$ , which implies that  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$  for all  $\mathbf{d}$  tangent to  $\Omega$  at  $\mathbf{x}^*$  (or, alternatively, for all  $\mathbf{d}$  orthogonal to  $\mathbf{x}^*$ ). ■

We now derive a second-order necessary condition that is satisfied by a local minimizer.

**Theorem 6.2 Second-Order Necessary Condition (SONC).** Let  $\Omega \subset \mathbb{R}^n$ ,  $f \in \mathcal{C}^2$  a function on  $\Omega$ ,  $\mathbf{x}^*$  a local minimizer of  $f$  over  $\Omega$ , and  $\mathbf{d}$  a feasible direction at  $\mathbf{x}^*$ . If  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$ , then

$$\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0,$$

where  $\mathbf{F}$  is the Hessian of  $f$ . □

*Proof.* We prove the result by contradiction. Suppose that there is a feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$  such that  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$  and  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} < 0$ . Let  $\mathbf{x}(\alpha) = \mathbf{x}^* + \alpha \mathbf{d}$  and define the composite function  $\phi(\alpha) = f(\mathbf{x}^* + \alpha \mathbf{d}) = f(\mathbf{x}(\alpha))$ . Then, by Taylor's theorem,

$$\phi(\alpha) = \phi(0) + \phi''(0) \frac{\alpha^2}{2} + o(\alpha^2),$$

where by assumption,  $\phi'(0) = \mathbf{d}^\top \nabla f(\mathbf{x}^*) = 0$  and  $\phi''(0) = \mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} < 0$ . For sufficiently small  $\alpha$ ,

$$\phi(\alpha) - \phi(0) = \phi''(0) \frac{\alpha^2}{2} + o(\alpha^2) < 0,$$

that is,

$$f(\mathbf{x}^* + \alpha \mathbf{d}) < f(\mathbf{x}^*),$$

which contradicts the assumption that  $\mathbf{x}^*$  is a local minimizer. Thus,

$$\phi''(0) = \mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0.$$
■

**Corollary 6.2 Interior Case.** Let  $\mathbf{x}^*$  be an interior point of  $\Omega \subset \mathbb{R}^n$ . If  $\mathbf{x}^*$  is a local minimizer of  $f : \Omega \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^2$ , then

$$\nabla f(\mathbf{x}^*) = \mathbf{0},$$

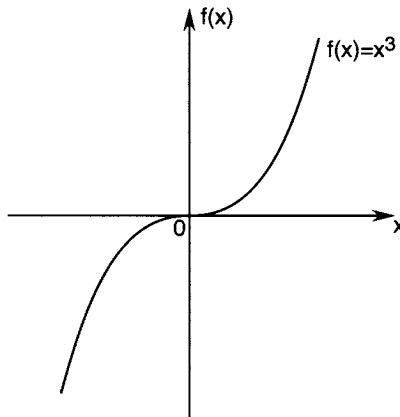
and  $\mathbf{F}(\mathbf{x}^*)$  is positive semidefinite ( $\mathbf{F}(\mathbf{x}^*) \geq 0$ ); that is, for all  $\mathbf{d} \in \mathbb{R}^n$ ,

$$\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{d} \geq 0.$$
□

*Proof.* If  $\mathbf{x}^*$  is an interior point, then all directions are feasible. The result then follows from Corollary 6.1 and Theorem 6.2. ■

In the examples below, we show that the necessary conditions are *not* sufficient.

**Example 6.6** Consider a function of one variable  $f(x) = x^3$ ,  $f : \mathbb{R} \rightarrow \mathbb{R}$ . Because  $f'(0) = 0$ , and  $f''(0) = 0$ , the point  $x = 0$  satisfies both the FONC and SONC. However,  $x = 0$  is not a minimizer (see Figure 6.6). ■

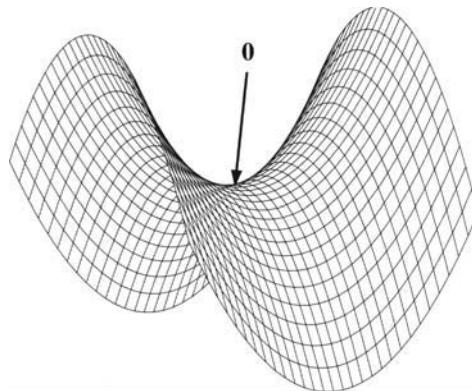


**Figure 6.6** The point  $0$  satisfies the FONC and SONC but is not a minimizer.

**Example 6.7** Consider a function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where  $f(\mathbf{x}) = x_1^2 - x_2^2$ . The FONC requires that  $\nabla f(\mathbf{x}) = [2x_1, -2x_2]^\top = \mathbf{0}$ . Thus,  $\mathbf{x} = [0, 0]^\top$  satisfies the FONC. The Hessian matrix of  $f$  is

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}.$$

The Hessian matrix is indefinite; that is, for some  $\mathbf{d}_1 \in \mathbb{R}^2$  we have  $\mathbf{d}_1^\top \mathbf{F} \mathbf{d}_1 > 0$  (e.g.,  $\mathbf{d}_1 = [1, 0]^\top$ ) and for some  $\mathbf{d}_2$  we have  $\mathbf{d}_2^\top \mathbf{F} \mathbf{d}_2 < 0$  (e.g.,  $\mathbf{d}_2 = [0, 1]^\top$ ). Thus,  $\mathbf{x} = [0, 0]^\top$  does not satisfy the SONC, and hence it is not a minimizer. The graph of  $f(\mathbf{x}) = x_1^2 - x_2^2$  is shown in Figure 6.7. ■



**Figure 6.7** Graph of  $f(\mathbf{x}) = x_1^2 - x_2^2$ . The point  $\mathbf{0}$  satisfies the FONC but not SONC; this point is not a minimizer.

We now derive sufficient conditions that imply that  $\mathbf{x}^*$  is a local minimizer.

**Theorem 6.3 Second-Order Sufficient Condition (SOSC), Interior Case.** Let  $f \in \mathcal{C}^2$  be defined on a region in which  $\mathbf{x}^*$  is an interior point. Suppose that

1.  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ .
2.  $\mathbf{F}(\mathbf{x}^*) > 0$ .

Then,  $\mathbf{x}^*$  is a strict local minimizer of  $f$ . □

*Proof.* Because  $f \in \mathcal{C}^2$ , we have  $\mathbf{F}(\mathbf{x}^*) = \mathbf{F}^\top(\mathbf{x}^*)$ . Using assumption 2 and Rayleigh's inequality it follows that if  $\mathbf{d} \neq \mathbf{0}$ , then  $0 < \lambda_{\min}(\mathbf{F}(\mathbf{x}^*))\|\mathbf{d}\|^2 \leq \mathbf{d}^\top \mathbf{F}(\mathbf{x}^*)\mathbf{d}$ . By Taylor's theorem and assumption 1,

$$f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*) = \frac{1}{2}\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*)\mathbf{d} + o(\|\mathbf{d}\|^2) \geq \frac{\lambda_{\min}(\mathbf{F}(\mathbf{x}^*))}{2}\|\mathbf{d}\|^2 + o(\|\mathbf{d}\|^2).$$

Hence, for all  $\mathbf{d}$  such that  $\|\mathbf{d}\|$  is sufficiently small,

$$f(\mathbf{x}^* + \mathbf{d}) > f(\mathbf{x}^*),$$

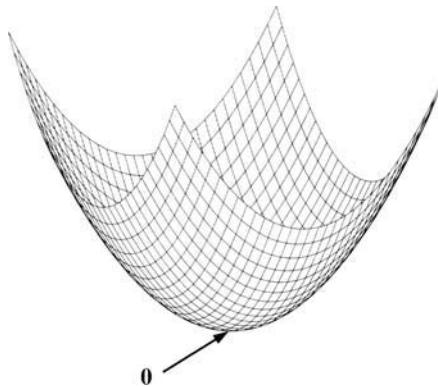
which completes the proof. ■

**Example 6.8** Let  $f(\mathbf{x}) = x_1^2 + x_2^2$ . We have  $\nabla f(\mathbf{x}) = [2x_1, 2x_2]^\top = \mathbf{0}$  if and only if  $\mathbf{x} = [0, 0]^\top$ . For all  $\mathbf{x} \in \mathbb{R}^2$ , we have

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} > 0.$$

The point  $\mathbf{x} = [0, 0]^\top$  satisfies the FONC, SONC, and SOSC. It is a strict local minimizer. Actually,  $\mathbf{x} = [0, 0]^\top$  is a strict global minimizer. Figure 6.8 shows the graph of  $f(\mathbf{x}) = x_1^2 + x_2^2$ . ■

In this chapter we presented a theoretical basis for the solution of nonlinear unconstrained problems. In the following chapters we are concerned with iterative methods of solving such problems. Such methods are of great importance in practice. Indeed, suppose that one is confronted with a highly nonlinear function of 20 variables. Then, the FONC requires the solution of 20 nonlinear simultaneous equations for 20 variables. These equations, being nonlinear, will normally have multiple solutions. In addition, we would have to compute 210 second derivatives (provided that  $f \in \mathcal{C}^2$ ) to use the SONC or SOSC. We begin our discussion of iterative methods in the next chapter with search methods for functions of one variable.



**Figure 6.8** Graph of  $f(\mathbf{x}) = x_1^2 + x_2^2$ .

## EXERCISES

**6.1** Consider the problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $f \in \mathcal{C}^2$ . For each of the following specifications for  $\Omega$ ,  $\mathbf{x}^*$ , and  $f$ , determine if the given point  $\mathbf{x}^*$  is: (i) definitely a local minimizer; (ii) definitely not a local minimizer; or (iii) possibly a local minimizer.

- a.  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 \geq 1\}$ ,  $\mathbf{x}^* = [1, 2]^\top$ , and gradient  $\nabla f(\mathbf{x}^*) = [1, 1]^\top$ .
- b.  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 \geq 1, x_2 \geq 2\}$ ,  $\mathbf{x}^* = [1, 2]^\top$ , and gradient  $\nabla f(\mathbf{x}^*) = [1, 0]^\top$ .
- c.  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 \geq 0, x_2 \geq 0\}$ ,  $\mathbf{x}^* = [1, 2]^\top$ , gradient  $\nabla f(\mathbf{x}^*) = [0, 0]^\top$ , and Hessian  $\mathbf{F}(\mathbf{x}^*) = \mathbf{I}$  (identity matrix).
- d.  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 \geq 1, x_2 \geq 2\}$ ,  $\mathbf{x}^* = [1, 2]^\top$ , gradient  $\nabla f(\mathbf{x}^*) = [1, 0]^\top$ , and Hessian

$$\mathbf{F}(\mathbf{x}^*) = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

**6.2** Find minimizers and maximizers of the function

$$f(x_1, x_2) = \frac{1}{3}x_1^3 - 4x_1 + \frac{1}{3}x_2^3 - 16x_2.$$

**6.3** Show that if  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ , and  $\mathbf{x}^* \in \Omega' \subset \Omega$ , then  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega'$ .

**6.4** Suppose that  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega$ , and  $\Omega \subset \Omega'$ . Show that if  $\mathbf{x}^*$  is an interior point of  $\Omega$ , then  $\mathbf{x}^*$  is a local minimizer of  $f$  over  $\Omega'$ . Show that the same conclusion cannot be made if  $\mathbf{x}^*$  is not an interior point of  $\Omega$ .

**6.5** Consider the problem of minimizing  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^3$ , over the constraint set  $\Omega$ . Suppose that 0 is an *interior point* of  $\Omega$ .

- a. Suppose that 0 is a local minimizer. By the FONC we know that  $f'(0) = 0$  (where  $f'$  is the first derivative of  $f$ ). By the SONC we know that  $f''(0) \geq 0$  (where  $f''$  is the second derivative of  $f$ ). State and prove a *third-order necessary condition (TONC)* involving the third derivative at 0,  $f'''(0)$ .
- b. Give an example of  $f$  such that the FONC, SONC, and TONC (in part a) hold at the interior point 0, but 0 is not a local minimizer of  $f$  over  $\Omega$ . (Show that your example is correct.)
- c. Suppose that  $f$  is a third-order polynomial. If 0 satisfies the FONC, SONC, and TONC (in part a), then is this *sufficient* for 0 to be a local minimizer?

**6.6** Consider the problem of minimizing  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^3$ , over the constraint set  $\Omega = [0, 1]$ . Suppose that  $x^* = 0$  is a local minimizer.

- a. By the FONC we know that  $f'(0) \geq 0$  (where  $f'$  is the first derivative of  $f$ ). By the SONC we know that if  $f'(0) = 0$ , then  $f''(0) \geq 0$  (where  $f''$  is the second derivative of  $f$ ). State and prove a *third-order necessary condition* involving the third derivative at 0,  $f'''(0)$ .
- b. Give an example of  $f$  such that the FONC, SONC, and TONC (in part a) hold at the point 0, but 0 is not a local minimizer of  $f$  over  $\Omega = [0, 1]$ .

**6.7** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{x}_0 \in \mathbb{R}^n$ , and  $\Omega \subset \mathbb{R}^n$ . Show that

$$\mathbf{x}_0 + \arg \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = \arg \min_{\mathbf{y} \in \Omega'} f(\mathbf{y}),$$

where  $\Omega' = \{\mathbf{y} : \mathbf{y} - \mathbf{x}_0 \in \Omega\}$ .

**6.8** Consider the following function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ :

$$f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 1 & 2 \\ 4 & 7 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 5 \end{bmatrix} + 6.$$

- a. Find the gradient and Hessian of  $f$  at the point  $[1, 1]^\top$ .
- b. Find the directional derivative of  $f$  at  $[1, 1]^\top$  with respect to a unit vector in the direction of maximal rate of increase.
- c. Find a point that satisfies the FONC (interior case) for  $f$ . Does this point satisfy the SONC (for a minimizer)?

**6.9** Consider the following function:

$$f(x_1, x_2) = x_1^2 x_2 + x_2^3 x_1.$$

- a. In what direction does the function  $f$  decrease most rapidly at the point  $\mathbf{x}^{(0)} = [2, 1]^\top$ ?
- b. What is the rate of increase of  $f$  at the point  $\mathbf{x}^{(0)}$  in the direction of maximum decrease of  $f$ ?
- c. Find the rate of increase of  $f$  at the point  $\mathbf{x}^{(0)}$  in the direction  $\mathbf{d} = [3, 4]^\top$ .

**6.10** Consider the following function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ :

$$f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 2 & 5 \\ -1 & 1 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 4 \end{bmatrix} + 7.$$

- a. Find the directional derivative of  $f$  at  $[0, 1]^\top$  in the direction  $[1, 0]^\top$ .
- b. Find all points that satisfy the first-order necessary condition for  $f$ . Does  $f$  have a minimizer? If it does, then find all minimizer(s); otherwise, explain why it does not.

**6.11** Consider the problem

$$\begin{aligned} &\text{minimize} && -x_2^2 \\ &\text{subject to} && |x_2| \leq x_1^2 \\ & && x_1 \geq 0, \end{aligned}$$

where  $x_1, x_2 \in \mathbb{R}$ .

- a. Does the point  $[x_1, x_2]^\top = \mathbf{0}$  satisfy the first-order necessary condition for a minimizer? That is, if  $f$  is the objective function, is it true that  $\mathbf{d}^\top \nabla f(\mathbf{0}) \geq 0$  for all feasible directions  $\mathbf{d}$  at  $\mathbf{0}$ ?
- b. Is the point  $[x_1, x_2]^\top = \mathbf{0}$  a local minimizer, a strict local minimizer, a local maximizer, a strict local maximizer, or none of the above?

**6.12** Consider the problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = 5x_2$  with  $\mathbf{x} = [x_1, x_2]^\top$ , and  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1^2 + x_2 \geq 1\}$ .

- a. Does the point  $\mathbf{x}^* = [0, 1]^\top$  satisfy the first-order necessary condition?
- b. Does the point  $\mathbf{x}^* = [0, 1]^\top$  satisfy the second-order necessary condition?
- c. Is the point  $\mathbf{x}^* = [0, 1]^\top$  a local minimizer?

**6.13** Consider the problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = -3x_1$  with  $\mathbf{x} = [x_1, x_2]^\top$ , and  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 + x_2^2 \leq 2\}$ . Answer each of the following questions, showing complete justification.

- a. Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the first-order necessary condition?
- b. Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the second-order necessary condition?
- c. Is the point  $\mathbf{x}^* = [2, 0]^\top$  a local minimizer?

**6.14** Consider the problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + x_2^2 \geq 1\}$  and  $f(\mathbf{x}) = x_2$ .

- a. Find all point(s) satisfying the FONC.
- b. Which of the point(s) in part a satisfy the SONC?
- c. Which of the point(s) in part a are local minimizers?

**6.15** Consider the problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega \end{aligned}$$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = 3x_1$  with  $\mathbf{x} = [x_1, x_2]^\top$ , and  $\Omega = \{\mathbf{x} = [x_1, x_2]^\top : x_1 + x_2^2 \geq 2\}$ . Answer each of the following questions, showing complete justification.

- Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the first-order necessary condition?
- Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the second-order necessary condition?
- Is the point  $\mathbf{x}^* = [2, 0]^\top$  a local minimizer?

*Hint:* Draw a picture with the constraint set and level sets of  $f$ .

**6.16** Consider the problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $\mathbf{x} = [x_1, x_2]^\top$ ,  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = 4x_1^2 - x_2^2$ , and  $\Omega = \{\mathbf{x} : x_1^2 + 2x_1 - x_2 \geq 0, x_1 \geq 0, x_2 \geq 0\}$ .

- Does the point  $\mathbf{x}^* = \mathbf{0} = [0, 0]^\top$  satisfy the first-order necessary condition?
- Does the point  $\mathbf{x}^* = \mathbf{0}$  satisfy the second-order necessary condition?
- Is the point  $\mathbf{x}^* = \mathbf{0}$  a local minimizer of the given problem?

**6.17** Consider the problem

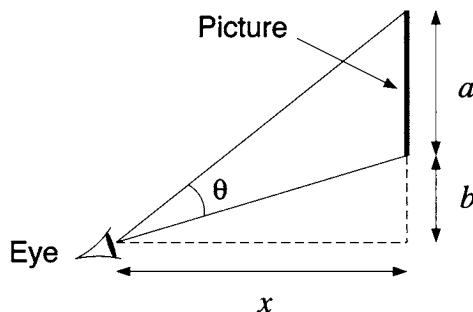
$$\begin{aligned} & \text{maximize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $\Omega \subset \{\mathbf{x} \in \mathbb{R}^2 : x_1 > 0, x_2 > 0\}$  and  $f : \Omega \rightarrow \mathbb{R}$  is given by  $f(\mathbf{x}) = \log(x_1) + \log(x_2)$  with  $\mathbf{x} = [x_1, x_2]^\top$ , where “log” represents natural logarithm. Suppose that  $\mathbf{x}^*$  is an optimal solution. Answer each of the following questions, showing complete justification.

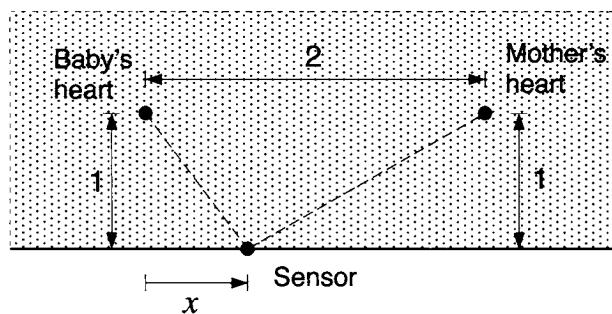
- Is it possible that  $\mathbf{x}^*$  is an interior point of  $\Omega$ ?
- At what point(s) (if any) is the second-order necessary condition satisfied?

**6.18** Suppose that we are given  $n$  real numbers,  $x_1, \dots, x_n$ . Find the number  $\bar{x} \in \mathbb{R}$  such that the sum of the squared difference between  $\bar{x}$  and the numbers above is minimized (assuming that the solution  $\bar{x}$  exists).

**6.19** An art collector stands at a distance of  $x$  feet from the wall, where a piece of art (picture) of height  $a$  feet is hung,  $b$  feet above his eyes, as shown in



**Figure 6.9** Art collector's eye position in Exercise 6.19.



**Figure 6.10** Simplified fetal heart monitoring system for Exercise 6.20.

Figure 6.9. Find the distance from the wall for which the angle  $\theta$  subtended by the eye to the picture is maximized.

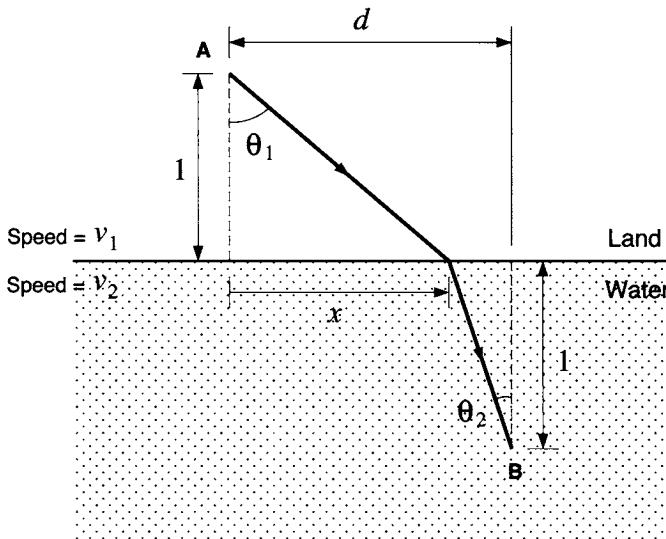
*Hint:* (1) Maximizing  $\theta$  is equivalent to maximizing  $\tan(\theta)$ .

(2) If  $\theta = \theta_2 - \theta_1$ , then  $\tan(\theta) = (\tan(\theta_2) - \tan(\theta_1))/(1 + \tan(\theta_2)\tan(\theta_1))$ .

**6.20** Figure 6.10 shows a simplified model of a fetal heart monitoring system (the distances shown have been scaled down to make the calculations simpler). A heartbeat sensor is located at position  $x$  (see Figure 6.10).

The energy of the heartbeat signal measured by the sensor is the reciprocal of the squared distance from the source (baby's heart or mother's heart). Find the position of the sensor that maximizes the *signal-to-interference ratio*, which is the ratio of the signal energy from the baby's heart to the signal energy from the mother's heart.

**6.21** An amphibian vehicle needs to travel from point A (on land) to point B (in water), as illustrated in Figure 6.11. The speeds at which the vehicle travels on land and water are  $v_1$  and  $v_2$ , respectively.



**Figure 6.11** Path of amphibian vehicle in Exercise 6.21.

- a. Suppose that the vehicle traverses a path that minimizes the total time taken to travel from A to B. Use the first-order necessary condition to show that for the optimal path above, the angles  $\theta_1$  and  $\theta_2$  in Figure 6.11 satisfy Snell's law:

$$\frac{\sin \theta_1}{\sin \theta_2} = \frac{v_1}{v_2}.$$

- b. Does the minimizer for the problem in part a satisfy the second-order sufficient condition?

**6.22** Suppose that you have a piece of land to sell and you have two buyers. If the first buyer receives a fraction  $x_1$  of the piece of land, the buyer will pay you  $U_1(x_1)$  dollars. Similarly, the second buyer will pay you  $U_2(x_2)$  dollars for a fraction of  $x_2$  of the land. Your goal is to sell parts of your land to the two buyers so that you maximize the total dollars you receive. (Other than the constraint that you can only sell whatever land you own, there are no restrictions on how much land you can sell to each buyer.)

- a. Formulate the problem as an optimization problem of the kind

$$\begin{aligned} & \text{maximize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega \end{aligned}$$

by specifying  $f$  and  $\Omega$ . Draw a picture of the constraint set.

- b. Suppose that  $U_i(x_i) = a_i x_i$ ,  $i = 1, 2$ , where  $a_1$  and  $a_2$  are given positive constants such that  $a_1 > a_2$ . Find all feasible points that satisfy the first-order necessary condition, giving full justification.
- c. Among those points in the answer of part b, find all that also satisfy the second-order necessary condition.

**6.23** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be defined by

$$f(\mathbf{x}) = (x_1 - x_2)^4 + x_1^2 - x_2^2 - 2x_1 + 2x_2 + 1,$$

where  $\mathbf{x} = [x_1, x_2]^\top$ . Suppose that we wish to minimize  $f$  over  $\mathbb{R}^2$ . Find all points satisfying the FONC. Do these points satisfy the SONC?

**6.24** Show that if  $\mathbf{d}$  is a feasible direction at a point  $\mathbf{x} \in \Omega$ , then for all  $\beta > 0$ , the vector  $\beta\mathbf{d}$  is also a feasible direction at  $\mathbf{x}$ .

**6.25** Let  $\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ . Show that  $\mathbf{d} \in \mathbb{R}^n$  is a feasible direction at  $\mathbf{x} \in \Omega$  if and only if  $\mathbf{A}\mathbf{d} = \mathbf{0}$ .

**6.26** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ . Consider the problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && x_1, x_2 \geq 0, \end{aligned}$$

where  $\mathbf{x} = [x_1, x_2]^\top$ . Suppose that  $\nabla f(\mathbf{0}) \neq \mathbf{0}$ , and

$$\frac{\partial f}{\partial x_1}(\mathbf{0}) \leq 0, \quad \frac{\partial f}{\partial x_2}(\mathbf{0}) \leq 0.$$

Show that  $\mathbf{0}$  cannot be a minimizer for this problem.

**6.27** Let  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{c} \neq \mathbf{0}$ , and consider the problem of minimizing the function  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$  over a constraint set  $\Omega \subset \mathbb{R}^n$ . Show that we cannot have a solution lying in the interior of  $\Omega$ .

**6.28** Consider the problem

$$\begin{aligned} &\text{maximize} && c_1 x_1 + c_2 x_2 \\ &\text{subject to} && x_1 + x_2 \leq 1 \\ & && x_1, x_2 \geq 0, \end{aligned}$$

where  $c_1$  and  $c_2$  are constants such that  $c_1 > c_2 \geq 0$ . This is a *linear programming* problem (see Part III). Assuming that the problem has an optimal feasible solution, use the first-order necessary condition to show that the *unique* optimal feasible solution  $\mathbf{x}^*$  is  $[1, 0]^\top$ .

*Hint:* First show that  $\mathbf{x}^*$  cannot lie in the interior of the constraint set. Then, show that  $\mathbf{x}^*$  cannot lie on the line segments  $L_1 = \{\mathbf{x} : x_1 = 0, 0 \leq x_2 < 1\}$ ,  $L_2 = \{\mathbf{x} : 0 \leq x_1 < 1, x_2 = 0\}$ ,  $L_3 = \{\mathbf{x} : 0 \leq x_1 < 1, x_2 = 1 - x_1\}$ .

**6.29 Line Fitting.** Let  $[x_1, y_1]^\top, \dots, [x_n, y_n]^\top$ ,  $n \geq 2$ , be points on the  $\mathbb{R}^2$  plane (each  $x_i, y_i \in \mathbb{R}$ ). We wish to find the straight line of “best fit” through these points (“best” in the sense that the average squared error is minimized); that is, we wish to find  $a, b \in \mathbb{R}$  to minimize

$$f(a, b) = \frac{1}{n} \sum_{i=1}^n (ax_i + b - y_i)^2.$$

a. Let

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i,$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

$$\bar{X^2} = \frac{1}{n} \sum_{i=1}^n x_i^2,$$

$$\bar{Y^2} = \frac{1}{n} \sum_{i=1}^n y_i^2,$$

$$\bar{XY} = \frac{1}{n} \sum_{i=1}^n x_i y_i.$$

Show that  $f(a, b)$  can be written in the form  $\mathbf{z}^\top \mathbf{Q} \mathbf{z} - 2\mathbf{c}^\top \mathbf{z} + d$ , where  $\mathbf{z} = [a, b]^\top$ ,  $\mathbf{Q} = \mathbf{Q}^\top \in \mathbb{R}^{2 \times 2}$ ,  $\mathbf{c} \in \mathbb{R}^2$  and  $d \in \mathbb{R}$ , and find expressions for  $\mathbf{Q}$ ,  $\mathbf{c}$ , and  $d$  in terms of  $\bar{X}$ ,  $\bar{Y}$ ,  $\bar{X^2}$ ,  $\bar{Y^2}$ , and  $\bar{XY}$ .

b. Assume that the  $x_i$ ,  $i = 1, \dots, n$ , are not all equal. Find the parameters  $a^*$  and  $b^*$  for the line of best fit in terms of  $\bar{X}$ ,  $\bar{Y}$ ,  $\bar{X^2}$ ,  $\bar{Y^2}$ , and  $\bar{XY}$ . Show that the point  $[a^*, b^*]^\top$  is the only local minimizer of  $f$ .

*Hint:*  $\bar{X^2} - (\bar{X})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$ .

c. Show that if  $a^*$  and  $b^*$  are the parameters of the line of best fit, then  $\bar{Y} = a^* \bar{X} + b^*$  (and hence once we have computed  $a^*$ , we can compute  $b^*$  using the formula  $b^* = \bar{Y} - a^* \bar{X}$ ).

**6.30** Suppose that we are given a set of vectors  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}\}$ ,  $\mathbf{x}^{(i)} \in \mathbb{R}^n$ ,  $i = 1, \dots, p$ . Find the vector  $\bar{\mathbf{x}} \in \mathbb{R}^n$  such that the average squared distance (norm) between  $\bar{\mathbf{x}}$  and  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}$ ,

$$\frac{1}{p} \sum_{i=1}^p \|\bar{\mathbf{x}} - \mathbf{x}^{(i)}\|^2,$$

is minimized. Use the SOSC to prove that the vector  $\bar{\mathbf{x}}$  found above is a strict local minimizer. How is  $\bar{\mathbf{x}}$  related to the centroid (or center of gravity) of the given set of points  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}\}$ ?

**6.31** Consider a function  $f : \Omega \rightarrow \mathbb{R}$ , where  $\Omega \subset \mathbb{R}^n$  is a convex set and  $f \in \mathcal{C}^1$ . Given  $\mathbf{x}^* \in \Omega$ , suppose that there exists  $c > 0$  such that  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq c\|\mathbf{d}\|$  for all feasible directions  $\mathbf{d}$  at  $\mathbf{x}^*$ . Show that  $\mathbf{x}^*$  is a strict local minimizer of  $f$  over  $\Omega$ .

**6.32** Prove the following generalization of the second-order sufficient condition:

**Theorem:** Let  $\Omega$  be a convex subset of  $\mathbb{R}^n$ ,  $f \in \mathcal{C}^2$  a real-valued function on  $\Omega$ , and  $\mathbf{x}^*$  a point in  $\Omega$ . Suppose that there exists  $c \in \mathbb{R}$ ,  $c > 0$ , such that for all feasible directions  $\mathbf{d}$  at  $\mathbf{x}^*$  ( $\mathbf{d} \neq \mathbf{0}$ ), the following hold:

1.  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0$ .
2.  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}^*)\mathbf{d} \geq c\|\mathbf{d}\|^2$ .

Then,  $\mathbf{x}^*$  is a strict local minimizer of  $f$ .

**6.33** Consider the quadratic function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Show that  $\mathbf{x}^*$  minimizes  $f$  if and only if  $\mathbf{x}^*$  satisfies the FONC.

**6.34** Consider the linear system  $x_{k+1} = ax_k + bu_{k+1}$ ,  $k \geq 0$ , where  $x_i \in \mathbb{R}$ ,  $u_i \in \mathbb{R}$ , and the initial condition is  $x_0 = 0$ . Find the values of the control inputs  $u_1, \dots, u_n$  to minimize

$$-qx_n + r \sum_{i=1}^n u_i^2,$$

where  $q, r > 0$  are given constants. This can be interpreted as desiring to make  $x_n$  as large as possible but at the same time desiring to make the total input energy  $\sum_{i=1}^n u_i^2$  as small as possible. The constants  $q$  and  $r$  reflect the relative weights of these two objectives.

# CHAPTER 7

---

## ONE-DIMENSIONAL SEARCH METHODS

---

### 7.1 Introduction

In this chapter, we are interested in the problem of minimizing an objective function  $f : \mathbb{R} \rightarrow \mathbb{R}$  (i.e., a one-dimensional problem). The approach is to use an iterative search algorithm, also called a line-search method. One-dimensional search methods are of interest for the following reasons. First, they are special cases of search methods used in multivariable problems. Second, they are used as part of general multivariable algorithms (as described later in Section 7.8).

In an iterative algorithm, we start with an initial candidate solution  $x^{(0)}$  and generate a sequence of *iterates*  $x^{(1)}, x^{(2)}, \dots$ . For each iteration  $k = 0, 1, 2, \dots$ , the next point  $x^{(k+1)}$  depends on  $x^{(k)}$  and the objective function  $f$ . The algorithm may use only the value of  $f$  at specific points, or perhaps its first derivative  $f'$ , or even its second derivative  $f''$ . In this chapter, we study several algorithms:

- Golden section method (uses only  $f$ )
- Fibonacci method (uses only  $f$ )

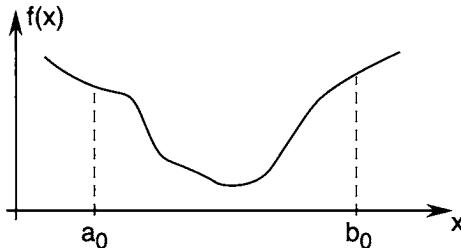


Figure 7.1 Unimodal function.

- Bisection method (uses only  $f'$ )
- Secant method (uses only  $f'$ )
- Newton's method (uses  $f'$  and  $f''$ )

The exposition here is based on [27].

## 7.2 Golden Section Search

The search methods we discuss in this and the next two sections allow us to determine the minimizer of an objective function  $f : \mathbb{R} \rightarrow \mathbb{R}$  over a closed interval, say  $[a_0, b_0]$ . The only property that we assume of the objective function  $f$  is that it is *unimodal*, which means that  $f$  has only one local minimizer. An example of such a function is depicted in Figure 7.1.

The methods we discuss are based on evaluating the objective function at different points in the interval  $[a_0, b_0]$ . We choose these points in such a way that an approximation to the minimizer of  $f$  may be achieved in as few evaluations as possible. Our goal is to narrow the range progressively until the minimizer is “boxed in” with sufficient accuracy.

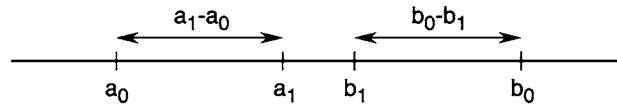
Consider a unimodal function  $f$  of one variable and the interval  $[a_0, b_0]$ . If we evaluate  $f$  at only one intermediate point of the interval, we cannot narrow the range within which we know the minimizer is located. We have to evaluate  $f$  at two intermediate points, as illustrated in Figure 7.2. We choose the intermediate points in such a way that the reduction in the range is symmetric, in the sense that

$$a_1 - a_0 = b_0 - b_1 = \rho(b_0 - a_0),$$

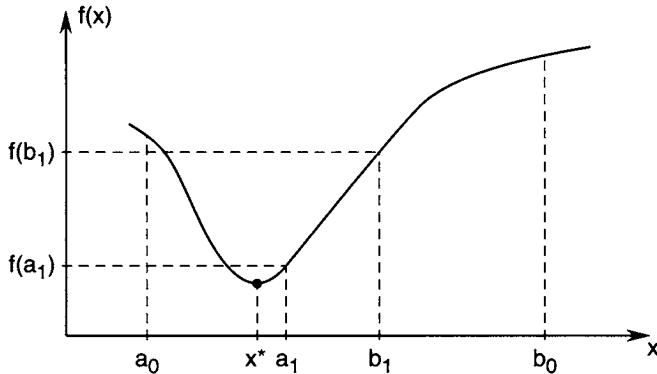
where

$$\rho < \frac{1}{2}.$$

We then evaluate  $f$  at the intermediate points. If  $f(a_1) < f(b_1)$ , then the minimizer must lie in the range  $[a_0, b_1]$ . If, on the other hand,  $f(a_1) \geq f(b_1)$ , then the minimizer is located in the range  $[a_1, b_0]$  (see Figure 7.3).



**Figure 7.2** Evaluating the objective function at two intermediate points.



**Figure 7.3** The case where  $f(a_1) < f(b_1)$ ; the minimizer  $x^* \in [a_0, b_1]$ .

Starting with the reduced range of uncertainty, we can repeat the process and similarly find two new points, say  $a_2$  and  $b_2$ , using the same value of  $\rho < \frac{1}{2}$  as before. However, we would like to minimize the number of objective function evaluations while reducing the width of the uncertainty interval. Suppose, for example, that  $f(a_1) < f(b_1)$ , as in Figure 7.3. Then, we know that  $x^* \in [a_0, b_1]$ . Because  $a_1$  is already in the uncertainty interval and  $f(a_1)$  is already known, we can make  $a_1$  coincide with  $b_2$ . Thus, only one new evaluation of  $f$  at  $a_2$  would be necessary. To find the value of  $\rho$  that results in only one new evaluation of  $f$ , see Figure 7.4. Without loss of generality, imagine that the original range  $[a_0, b_0]$  is of unit length. Then, to have only one new evaluation of  $f$  it is enough to choose  $\rho$  so that

$$\rho(b_1 - a_0) = b_1 - b_2.$$

Because  $b_1 - a_0 = 1 - \rho$  and  $b_1 - b_2 = 1 - 2\rho$ , we have

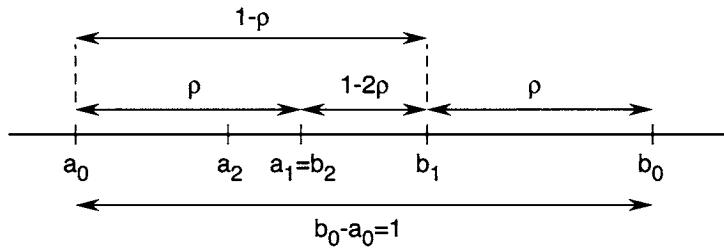
$$\rho(1 - \rho) = 1 - 2\rho.$$

We write the quadratic equation above as

$$\rho^2 - 3\rho + 1 = 0.$$

The solutions are

$$\rho_1 = \frac{3 + \sqrt{5}}{2}, \quad \rho_2 = \frac{3 - \sqrt{5}}{2}.$$



**Figure 7.4** Finding value of  $\rho$  resulting in only one new evaluation of  $f$ .

Because we require that  $\rho < \frac{1}{2}$ , we take

$$\rho = \frac{3 - \sqrt{5}}{2} \approx 0.382.$$

Observe that

$$1 - \rho = \frac{\sqrt{5} - 1}{2}$$

and

$$\frac{\rho}{1 - \rho} = \frac{3 - \sqrt{5}}{\sqrt{5} - 1} = \frac{\sqrt{5} - 1}{2} = \frac{1 - \rho}{1},$$

that is,

$$\frac{\rho}{1 - \rho} = \frac{1 - \rho}{1}.$$

Thus, dividing a range in the ratio of  $\rho$  to  $1 - \rho$  has the effect that the ratio of the shorter segment to the longer equals the ratio of the longer to the sum of the two. This rule was referred to by ancient Greek geometers as the *golden section*.

Using the golden section rule means that at every stage of the uncertainty range reduction (except the first), the objective function  $f$  need only be evaluated at one new point. The uncertainty range is reduced by the ratio  $1 - \rho \approx 0.61803$  at every stage. Hence,  $N$  steps of reduction using the golden section method reduces the range by the factor

$$(1 - \rho)^N \approx (0.61803)^N.$$

**Example 7.1** Suppose that we wish to use the golden section search method to find the value of  $x$  that minimizes

$$f(x) = x^4 - 14x^3 + 60x^2 - 70x$$

in the interval  $[0, 2]$  (this function comes from an example in [21]). We wish to locate this value of  $x$  to within a range of 0.3.

After  $N$  stages the range  $[0, 2]$  is reduced by  $(0.61803)^N$ . So, we choose  $N$  so that

$$(0.61803)^N \leq 0.3/2.$$

Four stages of reduction will do; that is,  $N = 4$ .

*Iteration 1.* We evaluate  $f$  at two intermediate points  $a_1$  and  $b_1$ . We have

$$\begin{aligned} a_1 &= a_0 + \rho(b_0 - a_0) = 0.7639, \\ b_1 &= a_0 + (1 - \rho)(b_0 - a_0) = 1.236, \end{aligned}$$

where  $\rho = (3 - \sqrt{5})/2$ . We compute

$$\begin{aligned} f(a_1) &= -24.36, \\ f(b_1) &= -18.96. \end{aligned}$$

Thus,  $f(a_1) < f(b_1)$ , so the uncertainty interval is reduced to

$$[a_0, b_1] = [0, 1.236].$$

*Iteration 2.* We choose  $b_2$  to coincide with  $a_1$ , and so  $f$  need only be evaluated at one new point,

$$a_2 = a_0 + \rho(b_1 - a_0) = 0.4721.$$

We have

$$\begin{aligned} f(a_2) &= -21.10, \\ f(b_2) &= f(a_1) = -24.36. \end{aligned}$$

Now,  $f(b_2) < f(a_2)$ , so the uncertainty interval is reduced to

$$[a_2, b_1] = [0.4721, 1.236].$$

*Iteration 3.* We set  $a_3 = b_2$  and compute  $b_3$ :

$$b_3 = a_2 + (1 - \rho)(b_1 - a_2) = 0.9443.$$

We have

$$\begin{aligned} f(a_3) &= f(b_2) = -24.36, \\ f(b_3) &= -23.59. \end{aligned}$$

So  $f(b_3) > f(a_3)$ . Hence, the uncertainty interval is further reduced to

$$[a_2, b_3] = [0.4721, 0.9443].$$

*Iteration 4.* We set  $b_4 = a_3$  and

$$a_4 = a_2 + \rho(b_3 - a_2) = 0.6525.$$

We have

$$\begin{aligned} f(a_4) &= -23.84, \\ f(b_4) = f(a_3) &= -24.36. \end{aligned}$$

Hence,  $f(a_4) > f(b_4)$ . Thus, the value of  $x$  that minimizes  $f$  is located in the interval

$$[a_4, b_3] = [0.6525, 0.9443].$$

Note that  $b_3 - a_4 = 0.292 < 0.3$ . ■

### 7.3 Fibonacci Method

Recall that the golden section method uses the same value of  $\rho$  throughout. Suppose now that we are allowed to vary the value  $\rho$  from stage to stage, so that at the  $k$ th stage in the reduction process we use a value  $\rho_k$ , at the next stage we use a value  $\rho_{k+1}$ , and so on.

As in the golden section search, our goal is to select successive values of  $\rho_k$ ,  $0 \leq \rho_k \leq 1/2$ , such that only one new function evaluation is required at each stage. To derive the strategy for selecting evaluation points, consider Figure 7.5. From this figure we see that it is sufficient to choose the  $\rho_k$  such that

$$\rho_{k+1}(1 - \rho_k) = 1 - 2\rho_k.$$

After some manipulations, we obtain

$$\rho_{k+1} = 1 - \frac{\rho_k}{1 - \rho_k}.$$

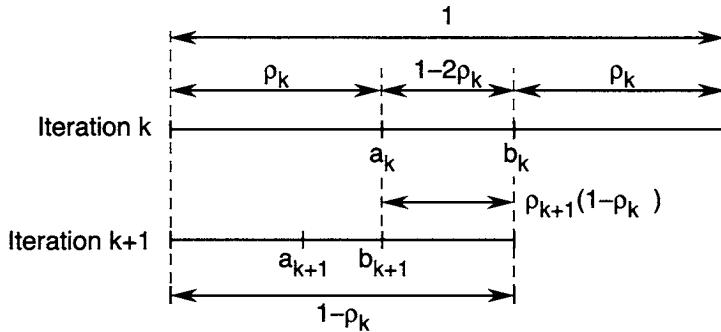
There are many sequences  $\rho_1, \rho_2, \dots$  that satisfy the law of formation above and the condition that  $0 \leq \rho_k \leq 1/2$ . For example, the sequence  $\rho_1 = \rho_2 = \rho_3 = \dots = (3 - \sqrt{5})/2$  satisfies the conditions above and gives rise to the golden section method.

Suppose that we are given a sequence  $\rho_1, \rho_2, \dots$  satisfying the conditions above and we use this sequence in our search algorithm. Then, after  $N$  iterations of the algorithm, the uncertainty range is reduced by a factor of

$$(1 - \rho_1)(1 - \rho_2) \cdots (1 - \rho_N).$$

Depending on the sequence  $\rho_1, \rho_2, \dots$ , we get a different reduction factor. The natural question is as follows: What sequence  $\rho_1, \rho_2, \dots$  minimizes the reduction factor above? This problem is a constrained optimization problem that can be stated formally as

$$\begin{aligned} &\text{minimize} && (1 - \rho_1)(1 - \rho_2) \cdots (1 - \rho_N) \\ &\text{subject to} && \rho_{k+1} = 1 - \frac{\rho_k}{1 - \rho_k}, \quad k = 1, \dots, N-1 \\ & && 0 \leq \rho_k \leq \frac{1}{2}, \quad k = 1, \dots, N. \end{aligned}$$



**Figure 7.5** Selecting evaluation points.

Before we give the solution to the optimization problem above, we need to introduce the *Fibonacci sequence*  $F_1, F_2, F_3, \dots$ . This sequence is defined as follows. First, let  $F_{-1} = 0$  and  $F_0 = 1$  by convention. Then, for  $k \geq 1$ ,

$$F_{k+1} = F_k + F_{k-1}.$$

Some values of elements in the Fibonacci sequence are:

$F_1$	$F_2$	$F_3$	$F_4$	$F_5$	$F_6$	$F_7$	$F_8$
1	2	3	5	8	13	21	34

It turns out that the solution to the optimization problem above is

$$\rho_1 = 1 - \frac{F_N}{F_{N+1}},$$

$$\rho_2 = 1 - \frac{F_{N-1}}{F_N},$$

⋮

$$\rho_k = 1 - \frac{F_{N-k+1}}{F_{N-k+2}},$$

⋮

$$\rho_N = 1 - \frac{F_1}{F_2},$$

where the  $F_k$  are the elements of the Fibonacci sequence. The resulting algorithm is called the *Fibonacci search method*. We present a proof for the optimality of the Fibonacci search method later in this section.

In the Fibonacci search method, the uncertainty range is reduced by the factor

$$(1 - \rho_1)(1 - \rho_2) \cdots (1 - \rho_N) = \frac{F_N}{F_{N+1}} \frac{F_{N-1}}{F_N} \cdots \frac{F_1}{F_2} = \frac{F_1}{F_{N+1}} = \frac{1}{F_{N+1}}.$$

Because the Fibonacci method uses the optimal values of  $\rho_1, \rho_2, \dots$ , the reduction factor above is less than that of the golden section method. In other words, the Fibonacci method is better than the golden section method in that it gives a smaller final uncertainty range.

We point out that there is an anomaly in the final iteration of the Fibonacci search method, because

$$\rho_N = 1 - \frac{F_1}{F_2} = \frac{1}{2}.$$

Recall that we need two intermediate points at each stage, one that comes from a previous iteration and another that is a new evaluation point. However, with  $\rho_N = 1/2$ , the two intermediate points coincide in the middle of the uncertainty interval, and therefore we cannot further reduce the uncertainty range. To get around this problem, we perform the new evaluation for the last iteration using  $\rho_N = 1/2 - \varepsilon$ , where  $\varepsilon$  is a small number. In other words, the new evaluation point is just to the left or right of the midpoint of the uncertainty interval. This modification to the Fibonacci method is, of course, of no significant practical consequence.

As a result of the modification above, the reduction in the uncertainty range at the last iteration may be either

$$1 - \rho_N = \frac{1}{2}$$

or

$$1 - (\rho_N - \varepsilon) = \frac{1}{2} + \varepsilon = \frac{1 + 2\varepsilon}{2},$$

depending on which of the two points has the smaller objective function value. Therefore, in the worst case, the reduction factor in the uncertainty range for the Fibonacci method is

$$\frac{1 + 2\varepsilon}{F_{N+1}}.$$

**Example 7.2** Consider the function

$$f(x) = x^4 - 14x^3 + 60x^2 - 70x.$$

Suppose that we wish to use the Fibonacci search method to find the value of  $x$  that minimizes  $f$  over the range  $[0, 2]$ , and locate this value of  $x$  to within the range 0.3.

After  $N$  steps the range is reduced by  $(1 + 2\varepsilon)/F_{N+1}$  in the worst case. We need to choose  $N$  such that

$$\frac{1 + 2\varepsilon}{F_{N+1}} \leq \frac{\text{final range}}{\text{initial range}} = \frac{0.3}{2} = 0.15.$$

Thus, we need

$$F_{N+1} \geq \frac{1 + 2\varepsilon}{0.15}.$$

If we choose  $\varepsilon \leq 0.1$ , then  $N = 4$  will do.

*Iteration 1.* We start with

$$1 - \rho_1 = \frac{F_4}{F_5} = \frac{5}{8}.$$

We then compute

$$\begin{aligned} a_1 &= a_0 + \rho_1(b_0 - a_0) = \frac{3}{4}, \\ b_1 &= a_0 + (1 - \rho_1)(b_0 - a_0) = \frac{5}{4}, \\ f(a_1) &= -24.34, \\ f(b_1) &= -18.65, \\ f(a_1) &< f(b_1). \end{aligned}$$

The range is reduced to

$$[a_0, b_1] = \left[0, \frac{5}{4}\right].$$

*Iteration 2.* We have

$$\begin{aligned} 1 - \rho_2 &= \frac{F_3}{F_4} = \frac{3}{5}, \\ a_2 &= a_0 + \rho_2(b_1 - a_0) = \frac{1}{2}, \\ b_2 &= a_1 = \frac{3}{4}, \\ f(a_2) &= -21.69, \\ f(b_2) &= f(a_1) = -24.34, \\ f(a_2) &> f(b_2), \end{aligned}$$

so the range is reduced to

$$[a_2, b_1] = \left[\frac{1}{2}, \frac{5}{4}\right].$$

*Iteration 3.* We compute

$$\begin{aligned} 1 - \rho_3 &= \frac{F_2}{F_3} = \frac{2}{3}, \\ a_3 &= b_2 = \frac{3}{4}, \\ b_3 &= a_2 + (1 - \rho_3)(b_1 - a_2) = 1, \\ f(a_3) &= f(b_2) = -24.34, \\ f(b_3) &= -23, \\ f(a_3) &< f(b_3). \end{aligned}$$

The range is reduced to

$$[a_2, b_3] = \left[ \frac{1}{2}, 1 \right].$$

*Iteration 4.* We choose  $\varepsilon = 0.05$ . We have

$$\begin{aligned} 1 - \rho_4 &= \frac{F_1}{F_2} = \frac{1}{2}, \\ a_4 &= a_2 + (\rho_4 - \varepsilon)(b_3 - a_2) = 0.725, \\ b_4 &= a_3 = \frac{3}{4}, \\ f(a_4) &= -24.27, \\ f(b_4) &= f(a_3) = -24.34, \\ f(a_4) &> f(b_4). \end{aligned}$$

The range is reduced to

$$[a_4, b_3] = [0.725, 1].$$

Note that  $b_3 - a_4 = 0.275 < 0.3$ . ■

We now turn to a proof of the optimality of the Fibonacci search method. Skipping the rest of this section does not affect the continuity of the presentation.

To begin, recall that we wish to prove that the values of  $\rho_1, \rho_2, \dots, \rho_N$  used in the Fibonacci method, where  $\rho_k = 1 - F_{N-k+1}/F_{N-k+2}$ , solve the optimization problem

$$\begin{aligned} &\text{minimize} \quad (1 - \rho_1)(1 - \rho_2) \cdots (1 - \rho_N) \\ &\text{subject to} \quad \rho_{k+1} = 1 - \frac{\rho_k}{1 - \rho_k}, \quad k = 1, \dots, N-1 \\ &\qquad \qquad \qquad 0 \leq \rho_k \leq \frac{1}{2}, \quad k = 1, \dots, N. \end{aligned}$$

It is easy to check that the values of  $\rho_1, \rho_2, \dots$  above for the Fibonacci search method satisfy the feasibility conditions in the optimization problem above (see Exercise 7.4). Recall that the Fibonacci method has an overall reduction factor of  $(1 - \rho_1) \cdots (1 - \rho_N) = 1/F_{N+1}$ . To prove that the Fibonacci search method is optimal, we show that for any feasible values of  $\rho_1, \dots, \rho_N$ , we have  $(1 - \rho_1) \cdots (1 - \rho_N) \geq 1/F_{N+1}$ .

It is more convenient to work with  $r_k = 1 - \rho_k$  rather than  $\rho_k$ . The optimization problem stated in terms of  $r_k$  is

$$\begin{aligned} &\text{minimize} \quad r_1 \cdots r_N \\ &\text{subject to} \quad r_{k+1} = \frac{1}{r_k} - 1, \quad k = 1, \dots, N-1 \\ &\qquad \qquad \qquad \frac{1}{2} \leq r_k \leq 1, \quad k = 1, \dots, N. \end{aligned}$$

Note that if  $r_1, r_2, \dots$  satisfy  $r_{k+1} = \frac{1}{r_k} - 1$ , then  $r_k \geq 1/2$  if and only if  $r_{k+1} \leq 1$ . Also,  $r_k \geq 1/2$  if and only if  $r_{k-1} \leq 2/3 \leq 1$ . Therefore, in the constraints above, we may remove the constraint  $r_k \leq 1$ , because it is implied implicitly by  $r_k \geq 1/2$  and the other constraints. Therefore, the constraints above reduce to

$$\begin{aligned} r_{k+1} &= \frac{1}{r_k} - 1, \quad k = 1, \dots, N-1, \\ r_k &\geq \frac{1}{2}, \quad k = 1, \dots, N. \end{aligned}$$

To proceed, we need the following technical lemmas. In the statements of the lemmas, we assume that  $r_1, r_2, \dots$  is a sequence that satisfies

$$r_{k+1} = \frac{1}{r_k} - 1, \quad r_k \geq \frac{1}{2}, \quad k = 1, 2, \dots.$$

**Lemma 7.1** *For  $k \geq 2$ ,*

$$r_k = -\frac{F_{k-2} - F_{k-1}r_1}{F_{k-3} - F_{k-2}r_1}.$$

□

*Proof.* We proceed by induction. For  $k = 2$  we have

$$r_2 = \frac{1}{r_1} - 1 = \frac{1 - r_1}{r_1} = -\frac{F_0 - F_1r_1}{F_{-1} - F_0r_1}$$

and hence the lemma holds for  $k = 2$ . Suppose now that the lemma holds for  $k \geq 2$ . We show that it also holds for  $k + 1$ . We have

$$\begin{aligned} r_{k+1} &= \frac{1}{r_k} - 1 \\ &= \frac{-F_{k-3} + F_{k-2}r_1}{F_{k-2} - F_{k-1}r_1} - \frac{F_{k-2} - F_{k-1}r_1}{F_{k-2} - F_{k-1}r_1} \\ &= -\frac{F_{k-2} + F_{k-3} - (F_{k-1} + F_{k-2})r_1}{F_{k-2} - F_{k-1}r_1} \\ &= -\frac{F_{k-1} - F_kr_1}{F_{k-2} - F_{k-1}r_1}, \end{aligned}$$

where we used the formation law for the Fibonacci sequence. ■

**Lemma 7.2** *For  $k \geq 2$ ,*

$$(-1)^k(F_{k-2} - F_{k-1}r_1) > 0.$$

□

*Proof.* We proceed by induction. For  $k = 2$ , we have

$$(-1)^2(F_0 - F_1 r_1) = 1 - r_1.$$

But  $r_1 = 1/(1 + r_2) \leq 2/3$ , and hence  $1 - r_1 > 0$ . Therefore, the result holds for  $k = 2$ . Suppose now that the lemma holds for  $k \geq 2$ . We show that it also holds for  $k + 1$ . We have

$$(-1)^{k+1}(F_{k-1} - F_k r_1) = (-1)^{k+1}r_{k+1} \frac{1}{r_{k+1}}(F_{k-1} - F_k r_1).$$

By Lemma 7.1,

$$r_{k+1} = -\frac{F_{k-1} - F_k r_1}{F_{k-2} - F_{k-1} r_1}.$$

Substituting for  $1/r_{k+1}$ , we obtain

$$(-1)^{k+1}(F_{k-1} - F_k r_1) = r_{k+1}(-1)^k(F_{k-2} - F_{k-1} r_1) > 0,$$

which completes the proof. ■

**Lemma 7.3** *For  $k \geq 2$ ,*

$$(-1)^{k+1}r_1 \geq (-1)^{k+1} \frac{F_k}{F_{k+1}}.$$

□

*Proof.* Because  $r_{k+1} = \frac{1}{r_k} - 1$  and  $r_k \geq \frac{1}{2}$ , we have  $r_{k+1} \leq 1$ . Substituting for  $r_{k+1}$  from Lemma 7.1, we get

$$-\frac{F_{k-1} - F_k r_1}{F_{k-2} - F_{k-1} r_1} \leq 1.$$

Multiplying the numerator and denominator by  $(-1)^k$  yields

$$\frac{(-1)^{k+1}(F_{k-1} - F_k r_1)}{(-1)^k(F_{k-2} - F_{k-1} r_1)} \leq 1.$$

By Lemma 7.2,  $(-1)^k(F_{k-2} - F_{k-1} r_1) > 0$ , and therefore we can multiply both sides of the inequality above by  $(-1)^k(F_{k-2} - F_{k-1} r_1)$  to obtain

$$(-1)^{k+1}(F_{k-1} - F_k r_1) \leq (-1)^k(F_{k-2} - F_{k-1} r_1).$$

Rearranging yields

$$(-1)^{k+1}(F_{k-1} + F_k)r_1 \geq (-1)^{k+1}(F_{k-2} + F_{k-1}).$$

Using the law of formation of the Fibonacci sequence, we get

$$(-1)^{k+1}F_{k+1}r_1 \geq (-1)^{k+1}F_k,$$

which upon dividing by  $F_{k+1}$  on both sides gives the desired result. ■

We are now ready to prove the optimality of the Fibonacci search method and the uniqueness of this optimal solution.

**Theorem 7.1** *Let  $r_1, \dots, r_N$ ,  $N \geq 2$ , satisfy the constraints*

$$\begin{aligned} r_{k+1} &= \frac{1}{r_k} - 1, \quad k = 1, \dots, N-1, \\ r_k &\geq \frac{1}{2}, \quad k = 1, \dots, N. \end{aligned}$$

*Then,*

$$r_1 \cdots r_N \geq \frac{1}{F_{N+1}}.$$

*Furthermore,*

$$r_1 \cdots r_N = \frac{1}{F_{N+1}}$$

*if and only if  $r_k = F_{N-k+1}/F_{N-k+2}$ ,  $k = 1, \dots, N$ . In other words, the values of  $r_1, \dots, r_N$  used in the Fibonacci search method form a unique solution to the optimization problem.* □

*Proof.* By substituting expressions for  $r_1, \dots, r_N$  from Lemma 7.1 and performing the appropriate cancellations, we obtain

$$r_1 \cdots r_N = (-1)^N (F_{N-2} - F_{N-1} r_1) = (-1)^N F_{N-2} + F_{N-1} (-1)^{N+1} r_1.$$

Using Lemma 7.3 yields

$$\begin{aligned} r_1 \cdots r_N &\geq (-1)^N F_{N-2} + F_{N-1} (-1)^{N+1} \frac{F_N}{F_{N+1}} \\ &= (-1)^N (F_{N-2} F_{N+1} - F_{N-1} F_N) \frac{1}{F_{N+1}}. \end{aligned}$$

By Exercise 7.5, it is readily checked that the following identity holds:  $(-1)^N (F_{N-2} F_{N+1} - F_{N-1} F_N) = 1$ . Hence,

$$r_1 \cdots r_N \geq \frac{1}{F_{N+1}}.$$

From the above we see that

$$r_1 \cdots r_N = \frac{1}{F_{N+1}}$$

if and only if

$$r_1 = \frac{F_N}{F_{N+1}}.$$

This is simply the value of  $r_1$  for the Fibonacci search method. Note that fixing  $r_1$  determines  $r_2, \dots, r_N$  uniquely. ■

For further discussion on the Fibonacci search method and its variants, see [133].

## 7.4 Bisection Method

Again we consider finding the minimizer of an objective function  $f : \mathbb{R} \rightarrow \mathbb{R}$  over an interval  $[a_0, b_0]$ . As before, we assume that the objective function  $f$  is unimodal. Further, suppose that  $f$  is continuously differentiable and that we can use values of the derivative  $f'$  as a basis for reducing the uncertainty interval.

The *bisection method* is a simple algorithm for successively reducing the uncertainty interval based on evaluations of the derivative. To begin, let  $x^{(0)} = (a_0 + b_0)/2$  be the midpoint of the initial uncertainty interval. Next, evaluate  $f'(x^{(0)})$ . If  $f'(x^{(0)}) > 0$ , then we deduce that the minimizer lies to the *left* of  $x^{(0)}$ . In other words, we reduce the uncertainty interval to  $[a_0, x^{(0)}]$ . On the other hand, if  $f'(x^{(0)}) < 0$ , then we deduce that the minimizer lies to the *right* of  $x^{(0)}$ . In this case, we reduce the uncertainty interval to  $[x^{(0)}, b_0]$ . Finally, if  $f'(x^{(0)}) = 0$ , then we declare  $x^{(0)}$  to be the minimizer and terminate our search.

With the new uncertainty interval computed, we repeat the process iteratively. At each iteration  $k$ , we compute the midpoint of the uncertainty interval. Call this point  $x^{(k)}$ . Depending on the sign of  $f'(x^{(k)})$  (assuming that it is nonzero), we reduce the uncertainty interval to the left or right of  $x^{(k)}$ . If at any iteration  $k$  we find that  $f'(x^{(k)}) = 0$ , then we declare  $x^{(k)}$  to be the minimizer and terminate our search.

Two salient features distinguish the bisection method from the golden section and Fibonacci methods. First, instead of using values of  $f$ , the bisection methods uses values of  $f'$ . Second, at each iteration, the length of the uncertainty interval is reduced by a factor of  $1/2$ . Hence, after  $N$  steps, the range is reduced by a factor of  $(1/2)^N$ . This factor is smaller than in the golden section and Fibonacci methods.

**Example 7.3** Recall Example 7.1 where we wish to find the minimizer of

$$f(x) = x^4 - 14x^3 + 60x^2 - 70x$$

in the interval  $[0, 2]$  to within a range of 0.3. The golden section method requires at least four stages of reduction. If, instead, we use the bisection method, we would choose  $N$  so that

$$(0.5)^N \leq 0.3/2.$$

In this case, only three stages of reduction are needed. ■

## 7.5 Newton's Method

Suppose again that we are confronted with the problem of minimizing a function  $f$  of a single real variable  $x$ . We assume now that at each measurement

point  $x^{(k)}$  we can determine  $f(x^{(k)})$ ,  $f'(x^{(k)})$ , and  $f''(x^{(k)})$ . We can fit a quadratic function through  $x^{(k)}$  that matches its first and second derivatives with that of the function  $f$ . This quadratic has the form

$$q(x) = f(x^{(k)}) + f'(x^{(k)})(x - x^{(k)}) + \frac{1}{2}f''(x^{(k)})(x - x^{(k)})^2.$$

Note that  $q(x^{(k)}) = f(x^{(k)})$ ,  $q'(x^{(k)}) = f'(x^{(k)})$ , and  $q''(x^{(k)}) = f''(x^{(k)})$ . Then, instead of minimizing  $f$ , we minimize its approximation  $q$ . The first-order necessary condition for a minimizer of  $q$  yields

$$0 = q'(x) = f'(x^{(k)}) + f''(x^{(k)})(x - x^{(k)}).$$

Setting  $x = x^{(k+1)}$ , we obtain

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

**Example 7.4** Using Newton's method, we will find the minimizer of

$$f(x) = \frac{1}{2}x^2 - \sin x.$$

Suppose that the initial value is  $x^{(0)} = 0.5$ , and that the required accuracy is  $\epsilon = 10^{-5}$ , in the sense that we stop when  $|x^{(k+1)} - x^{(k)}| < \epsilon$ .

We compute

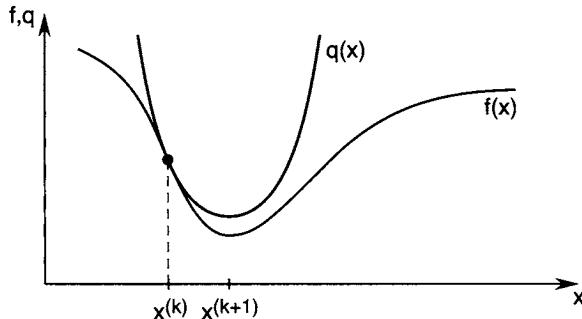
$$f'(x) = x - \cos x, \quad f''(x) = 1 + \sin x.$$

Hence,

$$\begin{aligned} x^{(1)} &= 0.5 - \frac{0.5 - \cos 0.5}{1 + \sin 0.5} \\ &= 0.5 - \frac{-0.3775}{1.479} \\ &= 0.7552. \end{aligned}$$

Proceeding in a similar manner, we obtain

$$\begin{aligned} x^{(2)} &= x^{(1)} - \frac{f'(x^{(1)})}{f''(x^{(1)})} = x^{(1)} - \frac{0.02710}{1.685} = 0.7391, \\ x^{(3)} &= x^{(2)} - \frac{f'(x^{(2)})}{f''(x^{(2)})} = x^{(2)} - \frac{9.461 \times 10^{-5}}{1.673} = 0.7390, \\ x^{(4)} &= x^{(3)} - \frac{f'(x^{(3)})}{f''(x^{(3)})} = x^{(3)} - \frac{1.17 \times 10^{-9}}{1.673} = 0.7390. \end{aligned}$$



**Figure 7.6** Newton's algorithm with  $f''(x) > 0$ .

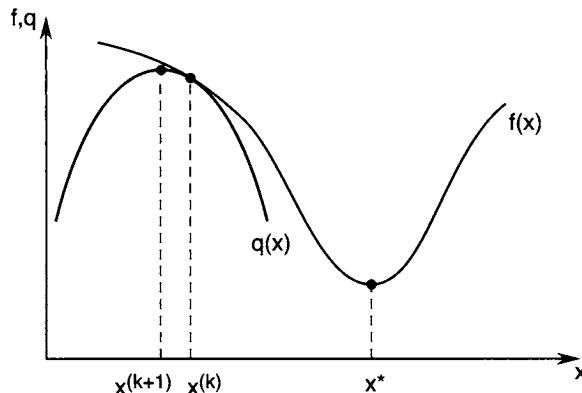
Note that  $|x^{(4)} - x^{(3)}| < \epsilon = 10^{-5}$ . Furthermore,  $f'(x^{(4)}) = -8.6 \times 10^{-6} \approx 0$ . Observe that  $f''(x^{(4)}) = 1.673 > 0$ , so we can assume that  $x^* \approx x^{(4)}$  is a strict minimizer. ■

Newton's method works well if  $f''(x) > 0$  everywhere (see Figure 7.6). However, if  $f''(x) < 0$  for some  $x$ , Newton's method may fail to converge to the minimizer (see Figure 7.7).

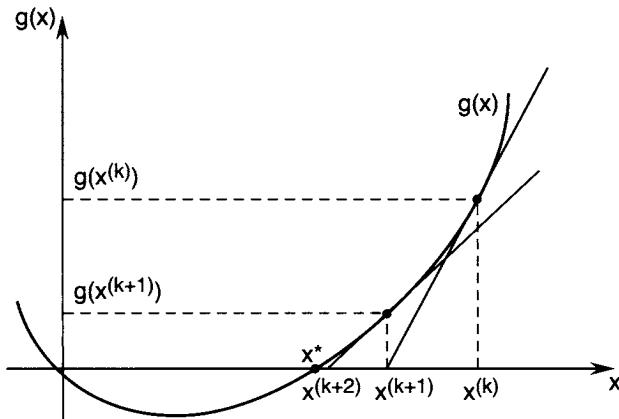
Newton's method can also be viewed as a way to drive the first derivative of  $f$  to zero. Indeed, if we set  $g(x) = f'(x)$ , then we obtain a formula for iterative solution of the equation  $g(x) = 0$ :

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})}.$$

In other words, we can use Newton's method for zero finding.



**Figure 7.7** Newton's algorithm with  $f''(x) < 0$ .



**Figure 7.8** Newton's method of tangents.

**Example 7.5** We apply Newton's method to improve a first approximation,  $x^{(0)} = 12$ , to the root of the equation

$$g(x) = x^3 - 12.2x^2 + 7.45x + 42 = 0.$$

We have  $g'(x) = 3x^2 - 24.4x + 7.45$ .

Performing two iterations yields

$$\begin{aligned} x^{(1)} &= 12 - \frac{102.6}{146.65} = 11.33, \\ x^{(2)} &= 11.33 - \frac{14.73}{116.11} = 11.21. \end{aligned}$$

■

Newton's method for solving equations of the form  $g(x) = 0$  is also referred to as *Newton's method of tangents*. This name is easily justified if we look at a geometric interpretation of the method when applied to the solution of the equation  $g(x) = 0$  (see Figure 7.8).

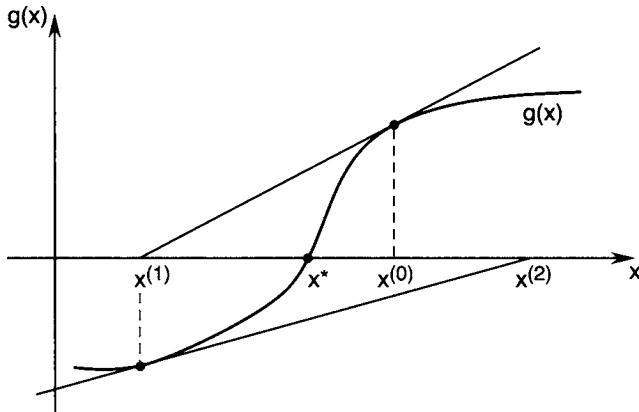
If we draw a tangent to  $g(x)$  at the given point  $x^{(k)}$ , then the tangent line intersects the  $x$ -axis at the point  $x^{(k+1)}$ , which we expect to be closer to the root  $x^*$  of  $g(x) = 0$ . Note that the slope of  $g(x)$  at  $x^{(k)}$  is

$$g'(x^{(k)}) = \frac{g(x^{(k)})}{x^{(k)} - x^{(k+1)}}.$$

Hence,

$$x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})}.$$

Newton's method of tangents may fail if the first approximation to the root is such that the ratio  $g(x^{(0)})/g'(x^{(0)})$  is not small enough (see Figure 7.9). Thus, an initial approximation to the root is very important.



**Figure 7.9** Example where Newton’s method of tangents fails to converge to the root  $x^*$  of  $g(x) = 0$ .

## 7.6 Secant Method

Newton’s method for minimizing  $f$  uses second derivatives of  $f$ :

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

If the second derivative is not available, we may attempt to approximate it using first derivative information. In particular, we may approximate  $f''(x^{(k)})$  above with

$$\frac{f'(x^{(k)}) - f'(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}.$$

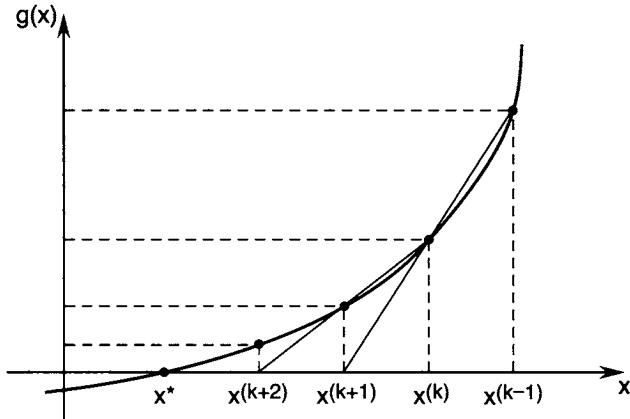
Using the foregoing approximation of the second derivative, we obtain the algorithm

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

called the *secant method*. Note that the algorithm requires two initial points to start it, which we denote  $x^{(-1)}$  and  $x^{(0)}$ . The secant algorithm can be represented in the following equivalent form:

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

Observe that, like Newton’s method, the secant method does not directly involve values of  $f(x^{(k)})$ . Instead, it tries to drive the derivative  $f'$  to zero. In fact, as we did for Newton’s method, we can interpret the secant method as an algorithm for solving equations of the form  $g(x) = 0$ . Specifically, the



**Figure 7.10** Secant method for root finding.

secant algorithm for finding a root of the equation  $g(x) = 0$  takes the form

$$x^{(k+1)} = x^{(k)} - \frac{x^{(k)} - x^{(k-1)}}{g(x^{(k)}) - g(x^{(k-1)})} g(x^{(k)}),$$

or, equivalently,

$$x^{(k+1)} = \frac{g(x^{(k)})x^{(k-1)} - g(x^{(k-1)})x^{(k)}}{g(x^{(k)}) - g(x^{(k-1)})}.$$

The secant method for root finding is illustrated in Figure 7.10 (compare this with Figure 7.8). Unlike Newton's method, which uses the slope of  $g$  to determine the next point, the secant method uses the “secant” between the  $(k - 1)$ th and  $k$ th points to determine the  $(k + 1)$ th point.

**Example 7.6** We apply the secant method to find the root of the equation

$$g(x) = x^3 - 12.2x^2 + 7.45x + 42 = 0.$$

We perform two iterations, with starting points  $x^{(-1)} = 13$  and  $x^{(0)} = 12$ . We obtain

$$\begin{aligned} x^{(1)} &= 11.40, \\ x^{(2)} &= 11.25. \end{aligned}$$

■

**Example 7.7** Suppose that the voltage across a resistor in a circuit decays according to the model  $V(t) = e^{-Rt}$ , where  $V(t)$  is the voltage at time  $t$  and  $R$  is the resistance value.

Given measurements  $V_1, \dots, V_n$  of the voltage at times  $t_1, \dots, t_n$ , respectively, we wish to find the best estimate of  $R$ . By the *best estimate* we mean the value of  $R$  that minimizes the total squared error between the measured voltages and the voltages predicted by the model.

We derive an algorithm to find the best estimate of  $R$  using the secant method. The objective function is

$$f(R) = \sum_{i=1}^n (V_i - e^{-Rt_i})^2.$$

Hence, we have

$$f'(R) = 2 \sum_{i=1}^n (V_i - e^{-Rt_i}) e^{-Rt_i} t_i.$$

The secant algorithm for the problem is

$$\begin{aligned} R_{k+1} &= R_k - \frac{R_k - R_{k-1}}{\sum_{i=1}^n (V_i - e^{-R_k t_i}) e^{-R_k t_i} t_i - (V_i - e^{-R_{k-1} t_i}) e^{-R_{k-1} t_i} t_i} \\ &\quad \times \sum_{i=1}^n (V_i - e^{-R_k t_i}) e^{-R_k t_i} t_i. \end{aligned}$$

■

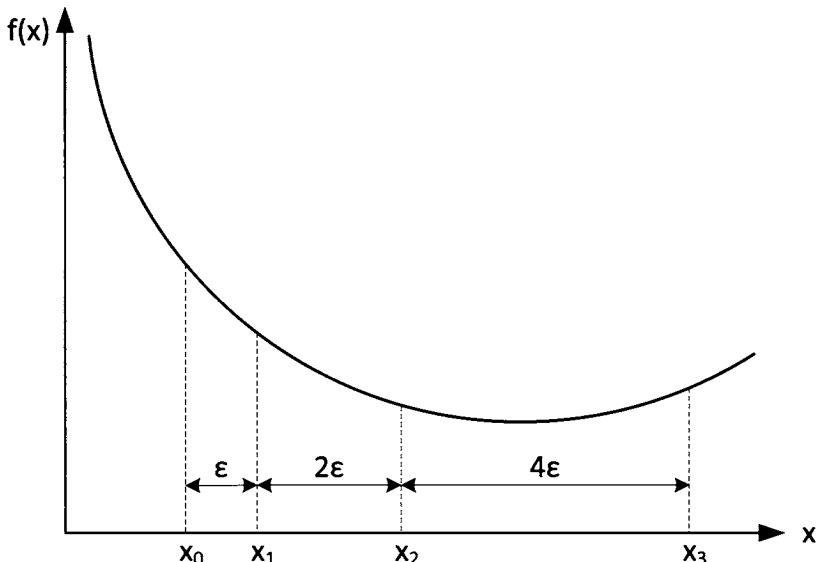
For further reading on the secant method, see [32]. Newton's method and the secant method are instances of *quadratic fit* methods. In Newton's method,  $x^{(k+1)}$  is the stationary point of a quadratic function that fits  $f'$  and  $f''$  at  $x^{(k)}$ . In the secant method,  $x^{(k+1)}$  is the stationary point of a quadratic function that fits  $f'$  at  $x^{(k)}$  and  $x^{(k-1)}$ . The secant method uses only  $f'$  (and not  $f''$ ) but needs values from two previous points. We leave it to the reader to verify that if we set  $x^{(k+1)}$  to be the stationary point of a quadratic function that fits  $f$  at  $x^{(k)}$ ,  $x^{(k-1)}$ , and  $x^{(k-2)}$ , we obtain a quadratic fit method that uses only values of  $f$ :

$$x^{(k+1)} = \frac{\sigma_{12} f(x^{(k)}) + \sigma_{20} f(x^{(k-1)}) + \sigma_{01} f(x^{(k-2)})}{2(\delta_{12} f(x^{(k)}) + \delta_{20} f(x^{(k-1)}) + \delta_{01} f(x^{(k-2)}))}$$

where  $\sigma_{ij} = (x^{(k-i)})^2 - (x^{(k-j)})^2$  and  $\delta_{ij} = x^{(k-i)} - x^{(k-j)}$  (see Exercise 7.9). This method does not use  $f'$  or  $f''$ , but needs values of  $f$  from three previous points. Three points are needed to initialize the iterations. The method is also sometimes called *inverse parabolic interpolation*.

An approach similar to fitting (or interpolation) based on higher-order polynomials is possible. For example, we could set  $x^{(k+1)}$  to be a stationary point of a *cubic* function that fits  $f'$  at  $x^{(k)}$ ,  $x^{(k-1)}$ , and  $x^{(k-2)}$ .

It is often practically advantageous to combine multiple methods, to overcome the limitations in any one method. For example, the golden section method is more robust but slower than inverse parabolic interpolation. *Brent's method* combines the two [17], resulting in a method that is faster than the golden section method but still retains its robustness properties.



**Figure 7.11** An illustration of the process of bracketing a minimizer.

## 7.7 Bracketing

Many of the methods we have described rely on an initial interval in which the minimizer is known to lie. This interval is also called a *bracket*, and procedures for finding such a bracket are called *bracketing* methods.

To find a bracket  $[a, b]$  containing the minimizer, assuming unimodality, it suffices to find three points  $a < c < b$  such that  $f(c) < f(a)$  and  $f(c) < f(b)$ . A simple bracketing procedure is as follows. First, we pick three arbitrary points  $x_0 < x_1 < x_2$ . If  $f(x_1) < f(x_0)$  and  $f(x_1) < f(x_2)$ , then we are done—the desired bracket is  $[x_0, x_2]$ . If not, say  $f(x_0) > f(x_1) > f(x_2)$ , then we pick a point  $x_3 > x_2$  and check if  $f(x_2) < f(x_3)$ . If it holds, then again we are done—the desired bracket is  $[x_1, x_3]$ . Otherwise, we continue with this process until the function increases. Typically, each new point chosen involves an expansion in distance between successive test points. For example, we could double the distance between successive points, as illustrated in Figure 7.11. An analogous process applies if the initial three points are such that  $f(x_0) < f(x_1) < f(x_2)$ .

In the procedure described above, when the bracketing process terminates, we have three points  $x_{k-2}$ ,  $x_{k-1}$ , and  $x_k$  such that  $f(x_{k-1}) < f(x_{k-2})$  and  $f(x_{k-1}) < f(x_k)$ . The desired bracket is then  $[x_{k-2}, x_k]$ , which we can then use to initialize any of a number of search methods, including the golden section, Fibonacci, and bisection methods. Note that at this point, we have already evaluated  $f(x_{k-2})$ ,  $f(x_{k-1})$ , and  $f(x_k)$ . If function evaluations are expensive to obtain, it would help if the point  $x_{k-1}$  inside the bracket also

coincides with one of the points used in the search method. For example, if we intend to use the golden section method, then it would help if  $x_{k-1} - x_{k-2} = \rho(x_k - x_{k-2})$ , where  $\rho = (3 - \sqrt{5})/2$ . In this case,  $x_{k-1}$  would be one of the two points within the initial interval used in the golden section method. This is achieved if each successive point  $x_k$  is chosen such that  $x_k = x_{k-1} + (2 - \rho)(x_{k-1} - x_{k-2})$ . In this case, the expansion in the distance between successive points is a factor  $2 - \rho \approx 1.618$ , which is less than double.

## 7.8 Line Search in Multidimensional Optimization

One-dimensional search methods play an important role in multidimensional optimization problems. In particular, iterative algorithms for solving such optimization problems (to be discussed in the following chapters) typically involve a *line search* at every iteration. To be specific, let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function that we wish to minimize. Iterative algorithms for finding a minimizer of  $f$  are of the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

where  $\mathbf{x}^{(0)}$  is a given initial point and  $\alpha_k \geq 0$  is chosen to minimize

$$\phi_k(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

The vector  $\mathbf{d}^{(k)}$  is called the *search direction* and  $\alpha_k$  is called the *step size*. Figure 7.12 illustrates a line search within a multidimensional setting. Note that choice of  $\alpha_k$  involves a one-dimensional minimization. This choice ensures that under appropriate conditions,

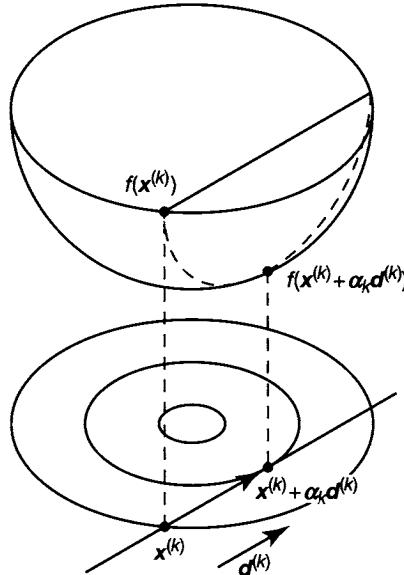
$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}).$$

Any of the one-dimensional methods discussed in this chapter (including bracketing) can be used to minimize  $\phi_k$ . We may, for example, use the secant method to find  $\alpha_k$ . In this case we need the derivative of  $\phi_k$ , which is

$$\phi'_k(\alpha) = \mathbf{d}^{(k)\top} \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

This is obtained using the chain rule. Therefore, applying the secant method for the line search requires the gradient  $\nabla f$ , the initial line-search point  $\mathbf{x}^{(k)}$ , and the search direction  $\mathbf{d}^{(k)}$  (see Exercise 7.11). Of course, other one-dimensional search methods may be used for line search (see, e.g., [43] and [88]).

Line-search algorithms used in practice involve considerations that we have not yet discussed thus far. First, determining the value of  $\alpha_k$  that exactly minimizes  $\phi_k$  may be computationally demanding; even worse, the minimizer of  $\phi_k$  may not even exist. Second, practical experience suggests that it is better to allocate more computational time on iterating the multidimensional



**Figure 7.12** Line search in multidimensional optimization.

optimization algorithm rather than performing exact line searches. These considerations led to the development of conditions for terminating line-search algorithms that would result in low-accuracy line searches while still securing a sufficient decrease in the value of the  $f$  from one iteration to the next. The basic idea is that we have to ensure that the step size  $\alpha_k$  is not too small or too large.

Some commonly used termination conditions are as follows. First, let  $\varepsilon \in (0, 1)$ ,  $\gamma > 1$ , and  $\eta \in (\varepsilon, 1)$  be given constants. The *Armijo condition* ensures that  $\alpha_k$  is not too large by requiring that

$$\phi_k(\alpha_k) \leq \phi_k(0) + \varepsilon \alpha_k \phi'_k(0).$$

Further, it ensures that  $\alpha_k$  is not too small by requiring that

$$\phi_k(\gamma \alpha_k) \geq \phi_k(0) + \varepsilon \gamma \alpha_k \phi'_k(0).$$

The *Goldstein condition* differs from Armijo in the second inequality:

$$\phi_k(\alpha_k) \geq \phi_k(0) + \eta \alpha_k \phi'_k(0).$$

The first Armijo inequality together with the Goldstein condition are often jointly called the *Armijo-Goldstein condition*. The *Wolfe condition* differs from Goldstein in that it involves only  $\phi'_k$ :

$$\phi'_k(\alpha_k) \geq \eta \phi'_k(0).$$

A stronger variation of this is the *strong Wolfe condition*:

$$|\phi'_k(\alpha_k)| \leq \eta |\phi'_k(0)|.$$

A simple practical (inexact) line-search method is the *Armijo backtracking algorithm*, described as follows. We start with some candidate value for the step size  $\alpha_k$ . If this candidate value satisfies a prespecified termination condition (usually the first Armijo inequality), then we stop and use it as the step size. Otherwise, we iteratively *decrease* the value by multiplying it by some constant factor  $\tau \in (0, 1)$  (typically  $\tau = 0.5$ ) and re-check the termination condition. If  $\alpha^{(0)}$  is the initial candidate value, then after  $m$  iterations the value obtained is  $\alpha_k = \tau^m \alpha^{(0)}$ . The algorithm *backtracks* from the initial value until the termination condition holds. In other words, the algorithm produces a value for the step size of the form  $\alpha_k = \tau^m \alpha^{(0)}$  with  $m$  being the smallest value in  $\{0, 1, 2, \dots\}$  for which  $\alpha_k$  satisfies the termination condition.

For more information on practical line-search methods, we refer the reader to [43, pp. 26–40], [96, Sec. 10.5], [11, App. C], [49], and [50].<sup>1</sup>

## EXERCISES

**7.1** Suppose that we have a unimodal function over the interval  $[5, 8]$ . Give an example of a desired final uncertainty range where the golden section method requires at least four iterations, whereas the Fibonacci method requires only three. You may choose an arbitrarily small value of  $\varepsilon$  for the Fibonacci method.

**7.2** Let  $f(x) = x^2 + 4 \cos x$ ,  $x \in \mathbb{R}$ . We wish to find the minimizer  $x^*$  of  $f$  over the interval  $[1, 2]$ . (*Calculator users*: Note that in  $\cos x$ , the argument  $x$  is in radians.)

- a. Plot  $f(x)$  versus  $x$  over the interval  $[1, 2]$ .
- b. Use the golden section method to locate  $x^*$  to within an uncertainty of 0.2. Display all intermediate steps using a table:

Iteration $k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New uncertainty interval
1	?	?	?	?	[?, ?]
2	?	?	?	?	[?, ?]
⋮	⋮	⋮	⋮	⋮	⋮

- c. Repeat part b using the Fibonacci method, with  $\varepsilon = 0.05$ . Display all intermediate steps using a table:

<sup>1</sup>We thank Dennis M. Goodman for furnishing us with references [49] and [50].

Iteration $k$	$\rho_k$	$a_k$	$b_k$	$f(a_k)$	$f(b_k)$	New uncertainty interval
1	?	?	?	?	?	[?,?]
2	?	?	?	?	?	[?,?]
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

- d. Apply Newton's method, using the same number of iterations as in part b, with  $x^{(0)} = 1$ .

**7.3** Let  $f(x) = 8e^{1-x} + 7\log(x)$ , where “log” represents the natural logarithm function.

- a. Use MATLAB to plot  $f(x)$  versus  $x$  over the interval  $[1, 2]$ , and verify that  $f$  is unimodal over  $[1, 2]$ .
- b. Write a simple MATLAB program to implement the golden section method that locates the minimizer of  $f$  over  $[1, 2]$  to within an uncertainty of 0.23. Display all intermediate steps using a table as in Exercise 7.2.
- c. Repeat part b using the Fibonacci method, with  $\varepsilon = 0.05$ . Display all intermediate steps using a table as in Exercise 7.2.

**7.4** Suppose that  $\rho_1, \dots, \rho_N$  are the values used in the Fibonacci search method. Show that for each  $k = 1, \dots, N$ ,  $0 \leq \rho_k \leq 1/2$ , and for each  $k = 1, \dots, N-1$ ,

$$\rho_{k+1} = 1 - \frac{\rho_k}{1 - \rho_k}.$$

**7.5** Show that if  $F_0, F_1, \dots$  is the Fibonacci sequence, then for each  $k = 2, 3, \dots$ ,

$$F_{k-2}F_{k+1} - F_{k-1}F_k = (-1)^k.$$

**7.6** Show that the Fibonacci sequence can be calculated using the formula

$$F_n = \frac{1}{\sqrt{5}} \left( \left( \frac{1+\sqrt{5}}{2} \right)^{n+1} - \left( \frac{1-\sqrt{5}}{2} \right)^{n+1} \right).$$

**7.7** Suppose that we have an efficient way of calculating exponentials. Based on this, use Newton's method to devise a method to approximate  $\log(2)$  [where “log” is the natural logarithm function]. Use an initial point of  $x^{(0)} = 1$ , and perform two iterations.

**7.8** Consider the problem of finding the zero of  $g(x) = (e^x - 1)/(e^x + 1)$ ,  $x \in \mathbb{R}$ , where  $e^x$  is the exponential of  $x$ . (Note that 0 is the unique zero of  $g$ .)

- a. Write down the algorithm for Newton's method of tangents applied to this problem. Simplify using the identity  $\sinh x = (e^x - e^{-x})/2$ .
- b. Find an initial condition  $x^{(0)}$  such that the algorithm cycles [i.e.,  $x^{(0)} = x^{(2)} = x^{(4)} = \dots$ ]. You need not explicitly calculate the initial condition; it suffices to provide an equation that the initial condition must satisfy.  
*Hint:* Draw a graph of  $g$ .
- c. For what values of the initial condition does the algorithm converge?

**7.9** Derive a one-dimensional search (minimization) algorithm based on quadratic fit with only objective function values. Specifically, derive an algorithm that computes  $x^{(k+1)}$  based on  $x^{(k)}$ ,  $x^{(k-1)}$ ,  $x^{(k-2)}$ ,  $f(x^{(k)})$ ,  $f(x^{(k-1)})$ , and  $f(x^{(k-2)})$ .

*Hint:* To simplify, use the notation  $\sigma_{ij} = (x^{(k-i)})^2 - (x^{(k-j)})^2$  and  $\delta_{ij} = x^{(k-i)} - x^{(k-j)}$ . You might also find it useful to experiment with your algorithm by writing a MATLAB program. Note that three points are needed to initialize the algorithm.

**7.10** The objective of this exercise is to implement the secant method using MATLAB.

- a. Write a simple MATLAB program to implement the secant method to locate the root of the equation  $g(x) = 0$ . For the stopping criterion, use the condition  $|x^{(k+1)} - x^{(k)}| < |x^{(k)}|\varepsilon$ , where  $\varepsilon > 0$  is a given constant.
- b. Let  $g(x) = (2x - 1)^2 + 4(4 - 1024x)^4$ . Find the root of  $g(x) = 0$  using the secant method with  $x^{(-1)} = 0$ ,  $x^{(0)} = 1$ , and  $\varepsilon = 10^{-5}$ . Also determine the value of  $g$  at the solution obtained.

**7.11** Write a MATLAB function that implements a line search algorithm using the secant method. The arguments to this function are the name of the M-file for the gradient, the current point, and the search direction. For example, the function may be called `linesearch_secant` and be used by the function call `alpha=linesearch_secant('grad',x,d)`, where `grad.m` is the M-file containing the gradient, `x` is the starting line search point, `d` is the search direction, and `alpha` is the value returned by the function [which we use in the following chapters as the step size for iterative algorithms (see, e.g., Exercises 8.25 and 10.11)].

*Note:* In the solutions manual, we used the stopping criterion  $|\mathbf{d}^\top \nabla f(\mathbf{x} + \alpha\mathbf{d})| \leq \varepsilon |\mathbf{d}^\top \nabla f(\mathbf{x})|$ , where  $\varepsilon > 0$  is a prespecified number,  $\nabla f$  is the gradient,  $\mathbf{x}$  is the starting line search point, and  $\mathbf{d}$  is the search direction. The rationale for the stopping criterion above is that we want to reduce the directional derivative of  $f$  in the direction  $\mathbf{d}$  by the specified fraction  $\varepsilon$ . We used a value of  $\varepsilon = 10^{-4}$  and initial conditions of 0 and 0.001.

**7.12** Consider using a gradient algorithm to minimize the function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{x}$$

with the initial guess  $\mathbf{x}^{(0)} = [0.8, -0.25]^\top$ .

- a. To initialize the line search, apply the bracketing procedure in Figure 7.11 along the line starting at  $\mathbf{x}^{(0)}$  in the direction of the negative gradient. Use  $\varepsilon = 0.075$ .
- b. Apply the golden section method to reduce the width of the uncertainty region to 0.01. Organize the results of your computation in a table format similar to that of Exercise 7.2.
- c. Repeat the above using the Fibonacci method.



# CHAPTER 8

---

## GRADIENT METHODS

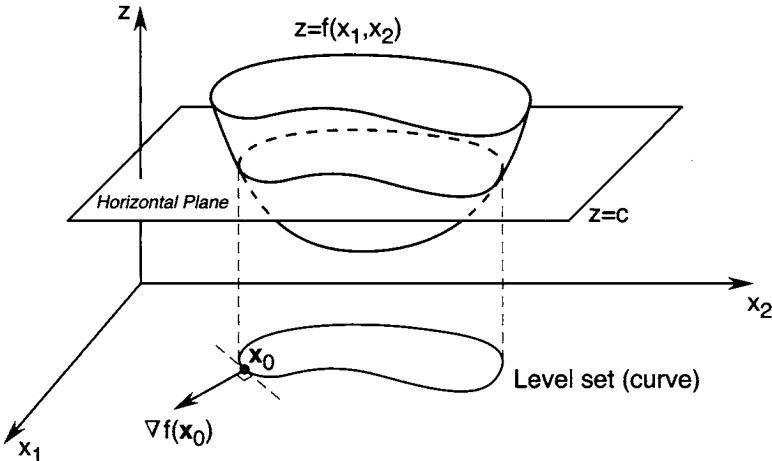
---

### 8.1 Introduction

In this chapter we consider a class of search methods for real-valued functions on  $\mathbb{R}^n$ . These methods use the gradient of the given function. In our discussion we use such terms as *level sets*, *normal vectors*, and *tangent vectors*. These notions were discussed in some detail in Part I.

Recall that a level set of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the set of points  $\mathbf{x}$  satisfying  $f(\mathbf{x}) = c$  for some constant  $c$ . Thus, a point  $\mathbf{x}_0 \in \mathbb{R}^n$  is on the level set corresponding to level  $c$  if  $f(\mathbf{x}_0) = c$ . In the case of functions of two real variables,  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , the notion of the level set is illustrated in Figure 8.1.

The gradient of  $f$  at  $\mathbf{x}_0$ , denoted  $\nabla f(\mathbf{x}_0)$ , if it is not a zero vector, is orthogonal to the tangent vector to an arbitrary smooth curve passing through  $\mathbf{x}_0$  on the level set  $f(\mathbf{x}) = c$ . Thus, the direction of maximum rate of increase of a real-valued differentiable function at a point is orthogonal to the level set of the function through that point. In other words, the gradient acts in such a direction that for a given small displacement, the function  $f$  increases more in the direction of the gradient than in any other direction. To prove this statement, recall that  $\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle, \|\mathbf{d}\| = 1$ , is the rate of increase of  $f$  in



**Figure 8.1** Constructing a level set corresponding to level  $c$  for  $f$ .

the direction  $\mathbf{d}$  at the point  $\mathbf{x}$ . By the Cauchy-Schwarz inequality,

$$\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle \leq \|\nabla f(\mathbf{x})\|$$

because  $\|\mathbf{d}\| = 1$ . But if  $\mathbf{d} = \nabla f(\mathbf{x})/\|\nabla f(\mathbf{x})\|$ , then

$$\left\langle \nabla f(\mathbf{x}), \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|} \right\rangle = \|\nabla f(\mathbf{x})\|.$$

Thus, the direction in which  $\nabla f(\mathbf{x})$  points is the direction of maximum rate of increase of  $f$  at  $\mathbf{x}$ . The direction in which  $-\nabla f(\mathbf{x})$  points is the direction of maximum rate of decrease of  $f$  at  $\mathbf{x}$ . Hence, the direction of negative gradient is a good direction to search if we want to find a function minimizer.

We proceed as follows. Let  $\mathbf{x}^{(0)}$  be a starting point, and consider the point  $\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$ . Then, by Taylor's theorem, we obtain

$$f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) = f(\mathbf{x}^{(0)}) - \alpha \|\nabla f(\mathbf{x}^{(0)})\|^2 + o(\alpha).$$

Thus, if  $\nabla f(\mathbf{x}^{(0)}) \neq \mathbf{0}$ , then for sufficiently small  $\alpha > 0$ , we have

$$f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) < f(\mathbf{x}^{(0)}).$$

This means that the point  $\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$  is an improvement over the point  $\mathbf{x}^{(0)}$  if we are searching for a minimizer.

To formulate an algorithm that implements this idea, suppose that we are given a point  $\mathbf{x}^{(k)}$ . To find the next point  $\mathbf{x}^{(k+1)}$ , we start at  $\mathbf{x}^{(k)}$  and move by an amount  $-\alpha_k \nabla f(\mathbf{x}^{(k)})$ , where  $\alpha_k$  is a positive scalar called the *step size*. This procedure leads to the following iterative algorithm:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}).$$

We refer to this as a *gradient descent algorithm* (or simply a *gradient algorithm*). The gradient varies as the search proceeds, tending to zero as we approach the minimizer. We have the option of either taking very small steps and reevaluating the gradient at every step, or we can take large steps each time. The first approach results in a laborious method of reaching the minimizer, whereas the second approach may result in a more zigzag path to the minimizer. The advantage of the second approach is possibly fewer gradient evaluations. Among many different methods that use this philosophy the most popular is the method of *steepest descent*, which we discuss next.

Gradient methods are simple to implement and often perform well. For this reason, they are used widely in practical applications. For a discussion of applications of the steepest descent method to the computation of optimal controllers, we recommend [85, pp. 481–515]. In Chapter 13 we apply a gradient method to the training of a class of neural networks.

## 8.2 The Method of Steepest Descent

The method of steepest descent is a gradient algorithm where the step size  $\alpha_k$  is chosen to achieve the maximum amount of decrease of the objective function at each individual step. Specifically,  $\alpha_k$  is chosen to minimize  $\phi_k(\alpha) \triangleq f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$ . In other words,

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)})).$$

To summarize, the steepest descent algorithm proceeds as follows: At each step, starting from the point  $\mathbf{x}^{(k)}$ , we conduct a line search in the direction  $-\nabla f(\mathbf{x}^{(k)})$  until a minimizer,  $\mathbf{x}^{(k+1)}$ , is found. A typical sequence resulting from the method of steepest descent is depicted in Figure 8.2.

Observe that the method of steepest descent moves in orthogonal steps, as stated in the following proposition.

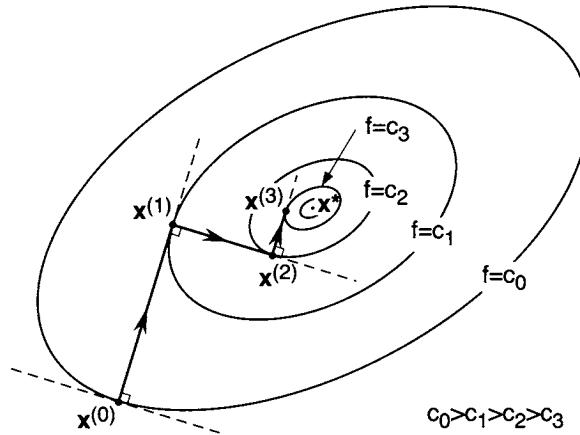
**Proposition 8.1** *If  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  is a steepest descent sequence for a given function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , then for each  $k$  the vector  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$  is orthogonal to the vector  $\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}$ .*  $\square$

*Proof.* From the iterative formula of the method of steepest descent it follows that

$$\langle \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)} \rangle = \alpha_k \alpha_{k+1} \langle \nabla f(\mathbf{x}^{(k)}), \nabla f(\mathbf{x}^{(k+1)}) \rangle.$$

To complete the proof it is enough to show that

$$\langle \nabla f(\mathbf{x}^{(k)}), \nabla f(\mathbf{x}^{(k+1)}) \rangle = 0.$$



**Figure 8.2** Typical sequence resulting from the method of steepest descent.

To this end, observe that  $\alpha_k$  is a nonnegative scalar that minimizes  $\phi_k(\alpha) \triangleq f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$ . Hence, using the FONC and the chain rule gives us

$$\begin{aligned} 0 &= \phi'_k(\alpha_k) \\ &= \frac{d\phi_k}{d\alpha}(\alpha_k) \\ &= \nabla f(\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}))^T (-\nabla f(\mathbf{x}^{(k)})) \\ &= -\langle \nabla f(\mathbf{x}^{(k+1)}), \nabla f(\mathbf{x}^{(k)}) \rangle, \end{aligned}$$

which completes the proof. ■

The proposition above implies that  $\nabla f(\mathbf{x}^{(k)})$  is parallel to the tangent plane to the level set  $\{f(\mathbf{x}) = f(\mathbf{x}^{(k+1)})\}$  at  $\mathbf{x}^{(k+1)}$ . Note that as each new point is generated by the steepest descent algorithm, the corresponding value of the function  $f$  decreases in value, as stated below.

**Proposition 8.2** *If  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  is the steepest descent sequence for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and if  $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , then  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ .* □

*Proof.* First recall that

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}),$$

where  $\alpha_k \geq 0$  is the minimizer of

$$\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$$

over all  $\alpha \geq 0$ . Thus, for  $\alpha \geq 0$ , we have

$$\phi_k(\alpha_k) \leq \phi_k(\alpha).$$

By the chain rule,

$$\phi'_k(0) = \frac{d\phi_k}{d\alpha}(0) = -(\nabla f(\mathbf{x}^{(k)}) - 0\nabla f(\mathbf{x}^{(k)}))^\top \nabla f(\mathbf{x}^{(k)}) = -\|\nabla f(\mathbf{x}^{(k)})\|^2 < 0$$

because  $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$  by assumption. Thus,  $\phi'_k(0) < 0$  and this implies that there is an  $\bar{\alpha} > 0$  such that  $\phi_k(0) > \phi_k(\alpha)$  for all  $\alpha \in (0, \bar{\alpha}]$ . Hence,

$$f(\mathbf{x}^{(k+1)}) = \phi_k(\alpha_k) \leq \phi_k(\bar{\alpha}) < \phi_k(0) = f(\mathbf{x}^{(k)}),$$

which completes the proof. ■

In Proposition 8.2, we proved that the algorithm possesses the *descent property*:  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$  if  $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ . If for some  $k$ , we have  $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$ , then the point  $\mathbf{x}^{(k)}$  satisfies the FONC. In this case,  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ . We can use the above as the basis for a stopping (termination) criterion for the algorithm.

The condition  $\nabla f(\mathbf{x}^{(k+1)}) = \mathbf{0}$ , however, is not directly suitable as a practical stopping criterion, because the numerical computation of the gradient will rarely be identically equal to zero. A practical stopping criterion is to check if the norm  $\|\nabla f(\mathbf{x}^{(k)})\|$  of the gradient is less than a prespecified threshold, in which case we stop. Alternatively, we may compute the absolute difference  $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|$  between objective function values for every two successive iterations, and if the difference is less than some prespecified threshold, then we stop; that is, we stop when

$$|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| < \varepsilon,$$

where  $\varepsilon > 0$  is a prespecified threshold. Yet another alternative is to compute the norm  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|$  of the difference between two successive iterates, and we stop if the norm is less than a prespecified threshold:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon.$$

Alternatively, we may check “relative” values of the quantities above; for example,

$$\frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{|f(\mathbf{x}^{(k)})|} < \varepsilon$$

or

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|} < \varepsilon.$$

The two (relative) stopping criteria above are preferable to the previous (absolute) criteria because the relative criteria are “scale-independent.” For example, scaling the objective function does not change the satisfaction of the criterion  $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|/|f(\mathbf{x}^{(k)})| < \varepsilon$ . Similarly, scaling the decision variable does not change the satisfaction of the criterion  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|/\|\mathbf{x}^{(k)}\| < \varepsilon$ .

To avoid dividing by very small numbers, we can modify these stopping criteria as follows:

$$\frac{|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|}{\max\{1, |f(\mathbf{x}^{(k)})|\}} < \varepsilon$$

or

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\max\{1, \|\mathbf{x}^{(k)}\|\}} < \varepsilon.$$

Note that the stopping criteria above are relevant to all the iterative algorithms we discuss in this part.

**Example 8.1** We use the method of steepest descent to find the minimizer of

$$f(x_1, x_2, x_3) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4.$$

The initial point is  $\mathbf{x}^{(0)} = [4, 2, -1]^\top$ . We perform three iterations.

We find that

$$\nabla f(\mathbf{x}) = [4(x_1 - 4)^3, 2(x_2 - 3), 16(x_3 + 5)^3]^\top.$$

Hence,

$$\nabla f(\mathbf{x}^{(0)}) = [0, -2, 1024]^\top.$$

To compute  $\mathbf{x}^{(1)}$ , we need

$$\begin{aligned}\alpha_0 &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) \\ &= \arg \min_{\alpha \geq 0} (0 + (2 + 2\alpha - 3)^2 + 4(-1 - 1024\alpha + 5)^4) \\ &= \arg \min_{\alpha \geq 0} \phi_0(\alpha).\end{aligned}$$

Using the secant method from Section 7.6, we obtain

$$\alpha_0 = 3.967 \times 10^{-3}.$$

For illustrative purpose, we show a plot of  $\phi_0(\alpha)$  versus  $\alpha$  in Figure 8.3, obtained using MATLAB. Thus,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}) = [4.000, 2.008, -5.062]^\top.$$

To find  $\mathbf{x}^{(2)}$ , we first determine

$$\nabla f(\mathbf{x}^{(1)}) = [0.000, -1.984, -0.003875]^\top.$$

Next, we find  $\alpha_1$ , where

$$\begin{aligned}\alpha_1 &= \arg \min_{\alpha \geq 0} (0 + (2.008 + 1.984\alpha - 3)^2 + 4(-5.062 + 0.003875\alpha + 5)^4) \\ &= \arg \min_{\alpha \geq 0} \phi_1(\alpha).\end{aligned}$$

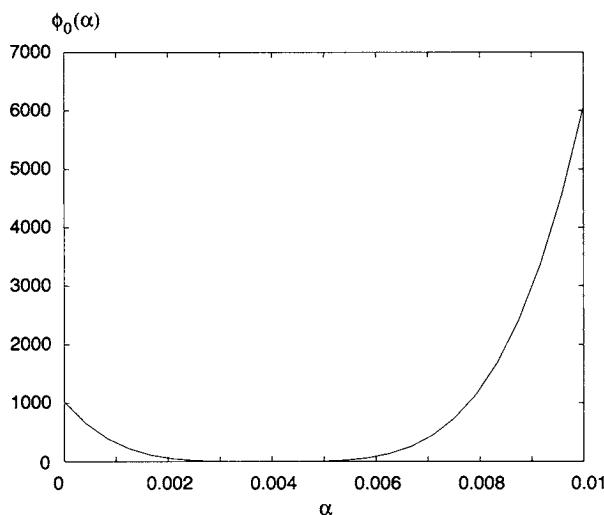


Figure 8.3 Plot of  $\phi_0(\alpha)$  versus  $\alpha$ .

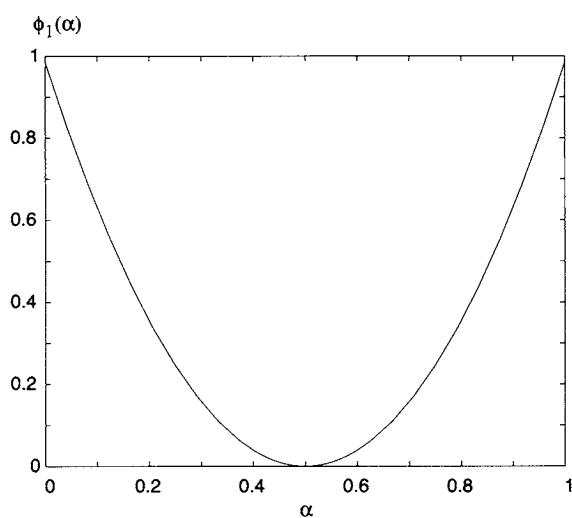
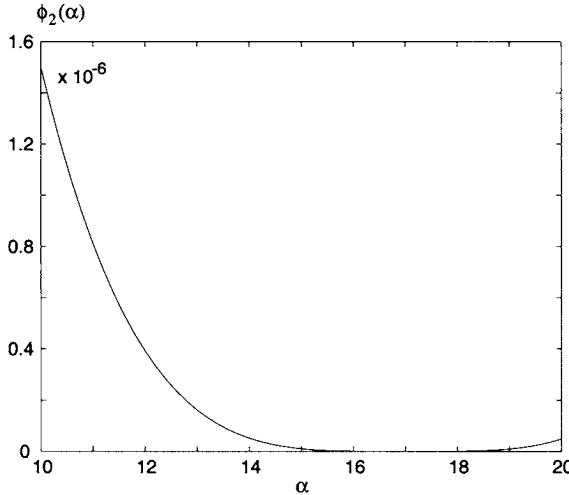


Figure 8.4 Plot of  $\phi_1(\alpha)$  versus  $\alpha$ .



**Figure 8.5** Plot of  $\phi_2(\alpha)$  versus  $\alpha$ .

Using the secant method again, we obtain  $\alpha_1 = 0.5000$ . Figure 8.4 depicts a plot of  $\phi_1(\alpha)$  versus  $\alpha$ . Thus,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha_1 \nabla f(\mathbf{x}^{(1)}) = [4.000, 3.000, -5.060]^\top.$$

To find  $\mathbf{x}^{(3)}$ , we first determine

$$\nabla f(\mathbf{x}^{(2)}) = [0.000, 0.000, -0.003525]^\top$$

and

$$\begin{aligned}\alpha_2 &= \arg \min_{\alpha \geq 0} (0.000 + 0.000 + 4(-5.060 + 0.003525\alpha + 5)^4) \\ &= \arg \min_{\alpha \geq 0} \phi_2(\alpha).\end{aligned}$$

We proceed as in the previous iterations to obtain  $\alpha_2 = 16.29$ . A plot of  $\phi_2(\alpha)$  versus  $\alpha$  is shown in Figure 8.5.

The value of  $\mathbf{x}^{(3)}$  is

$$\mathbf{x}^{(3)} = [4.000, 3.000, -5.002]^\top.$$

Note that the minimizer of  $f$  is  $[4, 3, -5]^\top$ , and hence it appears that we have arrived at the minimizer in only three iterations. The reader should be cautioned not to draw any conclusions from this example about the number of iterations required to arrive at a solution in general.

It goes without saying that numerical computations, such as those in this example, are performed in practice using a computer (rather than by hand).

The calculations above were written out explicitly, step by step, for the purpose of illustrating the operations involved in the steepest descent algorithm. The computations themselves were, in fact, carried out using a MATLAB program (see Exercise 8.25). ■

Let us now see what the method of steepest descent does with a quadratic function of the form

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is a symmetric positive definite matrix,  $\mathbf{b} \in \mathbb{R}^n$ , and  $\mathbf{x} \in \mathbb{R}^n$ . The unique minimizer of  $f$  can be found by setting the gradient of  $f$  to zero, where

$$\nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b},$$

because  $D(\mathbf{x}^\top \mathbf{Q} \mathbf{x}) = \mathbf{x}^\top (\mathbf{Q} + \mathbf{Q}^\top) = 2\mathbf{x}^\top \mathbf{Q}$ , and  $D(\mathbf{b}^\top \mathbf{x}) = \mathbf{b}^\top$ . There is no loss of generality in assuming  $\mathbf{Q}$  to be a symmetric matrix. For if we are given a quadratic form  $\mathbf{x}^\top \mathbf{A} \mathbf{x}$  and  $\mathbf{A} \neq \mathbf{A}^\top$ , then because the transposition of a scalar equals itself, we obtain

$$(\mathbf{x}^\top \mathbf{A} \mathbf{x})^\top = \mathbf{x}^\top \mathbf{A}^\top \mathbf{x} = \mathbf{x}^\top \mathbf{A} \mathbf{x}.$$

Hence,

$$\begin{aligned} \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{A}^\top \mathbf{x} \\ &= \frac{1}{2} \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \\ &\triangleq \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}. \end{aligned}$$

Note that

$$(\mathbf{A} + \mathbf{A}^\top)^\top = \mathbf{Q}^\top = \mathbf{A} + \mathbf{A}^\top = \mathbf{Q}.$$

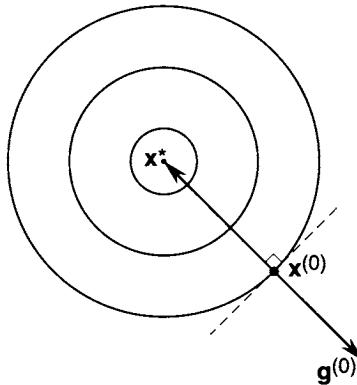
The Hessian of  $f$  is  $\mathbf{F}(\mathbf{x}) = \mathbf{Q} = \mathbf{Q}^\top > 0$ . To simplify the notation we write  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ . Then, the steepest descent algorithm for the quadratic function can be represented as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)},$$

where

$$\begin{aligned} \alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) \\ &= \arg \min_{\alpha \geq 0} \left( \frac{1}{2} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{Q} (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}) - (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{b} \right). \end{aligned}$$

In the quadratic case, we can find an explicit formula for  $\alpha_k$ . We proceed as follows. Assume that  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , for if  $\mathbf{g}^{(k)} = \mathbf{0}$ , then  $\mathbf{x}^{(k)} = \mathbf{x}^*$  and the



**Figure 8.6** Steepest descent method applied to  $f(x_1, x_2) = x_1^2 + x_2^2$ .

algorithm stops. Because  $\alpha_k \geq 0$  is a minimizer of  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})$ , we apply the FONC to  $\phi_k(\alpha)$  to obtain

$$\phi'_k(\alpha) = (\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})^\top \mathbf{Q}(-\mathbf{g}^{(k)}) - \mathbf{b}^\top (-\mathbf{g}^{(k)}).$$

Therefore,  $\phi'_k(\alpha) = 0$  if  $\alpha \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)} = (\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top) \mathbf{g}^{(k)}$ . But

$$\mathbf{x}^{(k)\top} \mathbf{Q} - \mathbf{b}^\top = \mathbf{g}^{(k)\top}.$$

Hence,

$$\alpha_k = \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}.$$

In summary, the method of steepest descent for the quadratic takes the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} \mathbf{g}^{(k)},$$

where

$$\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}.$$

**Example 8.2** Let

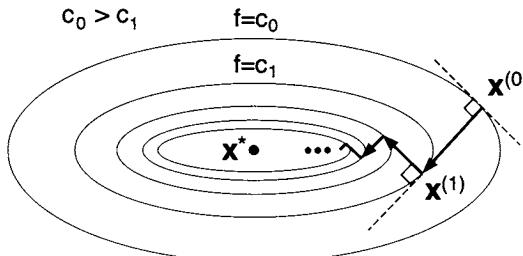
$$f(x_1, x_2) = x_1^2 + x_2^2.$$

Then, starting from an arbitrary initial point  $\mathbf{x}^{(0)} \in \mathbb{R}^2$ , we arrive at the solution  $\mathbf{x}^* = \mathbf{0} \in \mathbb{R}^2$  in only one step. See Figure 8.6.

However, if

$$f(x_1, x_2) = \frac{x_1^2}{5} + x_2^2,$$

then the method of steepest descent shuffles ineffectively back and forth when searching for the minimizer in a narrow valley (see Figure 8.7). This example illustrates a major drawback in the steepest descent method. More



**Figure 8.7** Steepest descent method in search for minimizer in a narrow valley.

sophisticated methods that alleviate this problem are discussed in subsequent chapters. ■

To understand better the method of steepest descent, we examine its convergence properties in the next section.

### 8.3 Analysis of Gradient Methods

#### Convergence

The method of steepest descent is an example of an iterative algorithm. This means that the algorithm generates a sequence of points, each calculated on the basis of the points preceding it. The method is a *descent method* because as each new point is generated by the algorithm, the corresponding value of the objective function decreases in value (i.e., the algorithm possesses the descent property).

We say that an iterative algorithm is *globally convergent* if for any arbitrary starting point the algorithm is guaranteed to generate a sequence of points converging to a point that satisfies the FONC for a minimizer. When the algorithm is not globally convergent, it may still generate a sequence that converges to a point satisfying the FONC, provided that the initial point is sufficiently close to the point. In this case we say that the algorithm is *locally convergent*. How close to a solution point we need to start for the algorithm to converge depends on the local convergence properties of the algorithm. A related issue of interest pertaining to a given locally or globally convergent algorithm is the *rate of convergence*; that is, how fast the algorithm converges to a solution point.

In this section we analyze the convergence properties of descent gradient methods, including the method of steepest descent and gradient methods with fixed step size. We can investigate important convergence characteristics of a gradient method by applying the method to quadratic problems. The convergence analysis is more convenient if instead of working with  $f$  we deal

with

$$V(\mathbf{x}) = f(\mathbf{x}) + \frac{1}{2} \mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^* = \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x} - \mathbf{x}^*),$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . The solution point  $\mathbf{x}^*$  is obtained by solving  $\mathbf{Q}\mathbf{x} = \mathbf{b}$ ; that is,  $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{b}$ . The function  $V$  differs from  $f$  only by a constant  $\frac{1}{2}\mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^*$ . We begin our analysis with the following useful lemma that applies to a general gradient algorithm.

**Lemma 8.1** *The iterative algorithm*

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$$

with  $\mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}$  satisfies

$$V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)}),$$

where if  $\mathbf{g}^{(k)} = \mathbf{0}$ , then  $\gamma_k = 1$ , and if  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , then

$$\gamma_k = \alpha_k \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)}} \left( 2 \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} - \alpha_k \right).$$

□

*Proof.* The proof is by direct computation. Note that if  $\mathbf{g}^{(k)} = \mathbf{0}$ , then the desired result holds trivially. In the remainder of the proof, assume that  $\mathbf{g}^{(k)} \neq \mathbf{0}$ . We first evaluate the expression

$$\frac{V(\mathbf{x}^{(k)}) - V(\mathbf{x}^{(k+1)})}{V(\mathbf{x}^{(k)})}.$$

To facilitate computations, let  $\mathbf{y}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^*$ . Then,  $V(\mathbf{x}^{(k)}) = \frac{1}{2} \mathbf{y}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)}$ . Hence,

$$\begin{aligned} V(\mathbf{x}^{(k+1)}) &= \frac{1}{2} (\mathbf{x}^{(k+1)} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x}^{(k+1)} - \mathbf{x}^*) \\ &= \frac{1}{2} (\mathbf{x}^{(k)} - \mathbf{x}^* - \alpha_k \mathbf{g}^{(k)})^\top \mathbf{Q} (\mathbf{x}^{(k)} - \mathbf{x}^* - \alpha_k \mathbf{g}^{(k)}) \\ &= \frac{1}{2} \mathbf{y}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)} - \alpha_k \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)} + \frac{1}{2} \alpha_k^2 \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}. \end{aligned}$$

Therefore,

$$\frac{V(\mathbf{x}^{(k)}) - V(\mathbf{x}^{(k+1)})}{V(\mathbf{x}^{(k)})} = \frac{2\alpha_k \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)} - \alpha_k^2 \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{y}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)}}.$$

Because

$$\mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{Q}\mathbf{x}^* = \mathbf{Q}\mathbf{y}^{(k)},$$

we have

$$\begin{aligned}\mathbf{y}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)} &= \mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)}, \\ \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{y}^{(k)} &= \mathbf{g}^{(k)\top} \mathbf{g}^{(k)}.\end{aligned}$$

Therefore, substituting the above, we get

$$\frac{V(\mathbf{x}^{(k)}) - V(\mathbf{x}^{(k+1)})}{V(\mathbf{x}^{(k)})} = \alpha_k \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)}} \left( 2 \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}} - \alpha_k \right) = \gamma_k.$$

■

Note that  $\gamma_k \leq 1$  for all  $k$ , because  $\gamma_k = 1 - V(\mathbf{x}^{(k+1)})/V(\mathbf{x}^{(k)})$  and  $V$  is a nonnegative function. If  $\gamma_k = 1$  for some  $k$ , then  $V(\mathbf{x}^{(k+1)}) = 0$ , which is equivalent to  $\mathbf{x}^{(k+1)} = \mathbf{x}^*$ . In this case we also have that for all  $i \geq k+1$ ,  $\mathbf{x}^{(i)} = \mathbf{x}^*$  and  $\gamma_i = 1$ . It turns out that  $\gamma_k = 1$  if and only if either  $\mathbf{g}^{(k)} = \mathbf{0}$  or  $\mathbf{g}^{(k)}$  is an eigenvector of  $\mathbf{Q}$  (see Lemma 8.3).

We are now ready to state and prove our key convergence theorem for gradient methods. The theorem gives a necessary and sufficient condition for the sequence  $\{\mathbf{x}^{(k)}\}$  generated by a gradient method to converge to  $\mathbf{x}^*$ ; that is,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  or  $\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*$ .

**Theorem 8.1** *Let  $\{\mathbf{x}^{(k)}\}$  be the sequence resulting from a gradient algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$ . Let  $\gamma_k$  be as defined in Lemma 8.1, and suppose that  $\gamma_k > 0$  for all  $k$ . Then,  $\{\mathbf{x}^{(k)}\}$  converges to  $\mathbf{x}^*$  for any initial condition  $\mathbf{x}^{(0)}$  if and only if*

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

□

*Proof.* From Lemma 8.1 we have  $V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)})$ , from which we obtain

$$V(\mathbf{x}^{(k)}) = \left( \prod_{i=0}^{k-1} (1 - \gamma_i) \right) V(\mathbf{x}^{(0)}).$$

Assume that  $\gamma_k < 1$  for all  $k$ , for otherwise the result holds trivially. Note that  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  if and only if  $V(\mathbf{x}^{(k)}) \rightarrow 0$ . By the equation above we see that this occurs if and only if  $\prod_{i=0}^{\infty} (1 - \gamma_i) = 0$ , which, in turn, holds if and only if  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$  (we get this simply by taking logs). Note that by Lemma 8.1,  $1 - \gamma_i \geq 0$  and  $\log(1 - \gamma_i)$  is well-defined [ $\log(0)$  is taken to be  $-\infty$ ]. Therefore, it remains to show that  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$  if and only if

$$\sum_{i=0}^{\infty} \gamma_i = \infty.$$

We first show that  $\sum_{i=0}^{\infty} \gamma_i = \infty$  implies that  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$ . For this, first observe that for any  $x \in \mathbb{R}$ ,  $x > 0$ , we have  $\log(x) \leq x - 1$

[this is easy to see simply by plotting  $\log(x)$  and  $x - 1$  versus  $x$ ]. Therefore,  $\log(1 - \gamma_i) \leq 1 - \gamma_i - 1 = -\gamma_i$ , and hence  $-\log(1 - \gamma_i) \geq \gamma_i$ . Thus, if  $\sum_{i=0}^{\infty} \gamma_i = \infty$ , then clearly  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$ .

Finally, we show that  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) = \infty$  implies that  $\sum_{i=0}^{\infty} \gamma_i = \infty$ . We proceed by contraposition. Suppose that  $\sum_{i=0}^{\infty} \gamma_i < \infty$ . Then, it must be that  $\gamma_i \rightarrow 0$ . Now observe that for  $x \in \mathbb{R}$ ,  $x \leq 1$  and  $x$  sufficiently close to 1, we have  $\log(x) \geq 2(x - 1)$  [as before, this is easy to see simply by plotting  $\log(x)$  and  $2(x - 1)$  versus  $x$ ]. Therefore, for sufficiently large  $i$ ,  $\log(1 - \gamma_i) \geq 2(1 - \gamma_i - 1) = -2\gamma_i$ , which implies that  $-\log(1 - \gamma_i) \leq 2\gamma_i$ . Hence,  $\sum_{i=0}^{\infty} \gamma_i < \infty$  implies that  $\sum_{i=0}^{\infty} -\log(1 - \gamma_i) < \infty$ .

This completes the proof. ■

The assumption in Theorem 8.1 that  $\gamma_k > 0$  for all  $k$  is significant in that it corresponds to the algorithm having the descent property (see Exercise 8.23). Furthermore, the result of the theorem does not hold in general if we do not assume that  $\gamma_k > 0$  for all  $k$ , as shown in the following example.

**Example 8.3** We show, using a counterexample, that the assumption that  $\gamma_k > 0$  in Theorem 8.1 is necessary for the result of the theorem to hold. Indeed, for each  $k = 0, 1, 2, \dots$ , choose  $\alpha_k$  in such a way that  $\gamma_{2k} = -1/2$  and  $\gamma_{2k+1} = 1/2$  (we can always do this if, for example,  $\mathbf{Q} = \mathbf{I}_n$ ). From Lemma 8.1 we have

$$V(\mathbf{x}^{(2(k+1))}) = (1 - 1/2)(1 + 1/2)V(\mathbf{x}^{(2k)}) = (3/4)V(\mathbf{x}^{(2k)}).$$

Therefore,  $V(\mathbf{x}^{(2k)}) \rightarrow 0$ . Because  $V(\mathbf{x}^{(2k+1)}) = (3/2)V(\mathbf{x}^{(2k)})$ , we also have that  $V(\mathbf{x}^{(2k+1)}) \rightarrow 0$ . Hence,  $V(\mathbf{x}^{(k)}) \rightarrow 0$ , which implies that  $\mathbf{x}^{(k)} \rightarrow 0$  (for all  $\mathbf{x}^{(0)}$ ). On the other hand, it is clear that

$$\sum_{i=0}^k \gamma_i \leq \frac{1}{2}$$

for all  $k$ . Hence, the result of the theorem does not hold if  $\gamma_k \leq 0$  for some  $k$ . ■

Using the general theorem above, we can now establish the convergence of specific cases of the gradient algorithm, including the steepest descent algorithm and algorithms with fixed step size. In the analysis to follow, we use Rayleigh's inequality, which states that for any  $\mathbf{Q} = \mathbf{Q}^\top > 0$ , we have

$$\lambda_{\min}(\mathbf{Q})\|\mathbf{x}\|^2 \leq \mathbf{x}^\top \mathbf{Q} \mathbf{x} \leq \lambda_{\max}(\mathbf{Q})\|\mathbf{x}\|^2,$$

where  $\lambda_{\min}(\mathbf{Q})$  denotes the minimal eigenvalue of  $\mathbf{Q}$  and  $\lambda_{\max}(\mathbf{Q})$  denotes the maximal eigenvalue of  $\mathbf{Q}$ . For  $\mathbf{Q} = \mathbf{Q}^\top > 0$ , we also have

$$\begin{aligned} \lambda_{\min}(\mathbf{Q}^{-1}) &= \frac{1}{\lambda_{\max}(\mathbf{Q})}, \\ \lambda_{\max}(\mathbf{Q}^{-1}) &= \frac{1}{\lambda_{\min}(\mathbf{Q})}, \end{aligned}$$

and

$$\lambda_{\min}(\mathbf{Q}^{-1})\|\mathbf{x}\|^2 \leq \mathbf{x}^\top \mathbf{Q}^{-1} \mathbf{x} \leq \lambda_{\max}(\mathbf{Q}^{-1})\|\mathbf{x}\|^2.$$

**Lemma 8.2** Let  $\mathbf{Q} = \mathbf{Q}^\top > 0$  be an  $n \times n$  real symmetric positive definite matrix. Then, for any  $\mathbf{x} \in \mathbb{R}^n$ , we have

$$\frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})} \leq \frac{(\mathbf{x}^\top \mathbf{x})^2}{(\mathbf{x}^\top \mathbf{Q} \mathbf{x})(\mathbf{x}^\top \mathbf{Q}^{-1} \mathbf{x})} \leq \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}.$$

□

*Proof.* Applying Rayleigh's inequality and using the properties of symmetric positive definite matrices listed previously, we get

$$\frac{(\mathbf{x}^\top \mathbf{x})^2}{(\mathbf{x}^\top \mathbf{Q} \mathbf{x})(\mathbf{x}^\top \mathbf{Q}^{-1} \mathbf{x})} \leq \frac{\|\mathbf{x}\|^4}{\lambda_{\min}(\mathbf{Q})\|\mathbf{x}\|^2\lambda_{\min}(\mathbf{Q}^{-1})\|\mathbf{x}\|^2} = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$$

and

$$\frac{(\mathbf{x}^\top \mathbf{x})^2}{(\mathbf{x}^\top \mathbf{Q} \mathbf{x})(\mathbf{x}^\top \mathbf{Q}^{-1} \mathbf{x})} \geq \frac{\|\mathbf{x}\|^4}{\lambda_{\max}(\mathbf{Q})\|\mathbf{x}\|^2\lambda_{\max}(\mathbf{Q}^{-1})\|\mathbf{x}\|^2} = \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})}.$$

■

We are now ready to establish the convergence of the steepest descent method.

**Theorem 8.2** In the steepest descent algorithm, we have  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for any  $\mathbf{x}^{(0)}$ . □

*Proof.* If  $\mathbf{g}^{(k)} = \mathbf{0}$  for some  $k$ , then  $\mathbf{x}^{(k)} = \mathbf{x}^*$  and the result holds. So assume that  $\mathbf{g}^{(k)} \neq \mathbf{0}$  for all  $k$ . Recall that for the steepest descent algorithm,

$$\alpha_k = \frac{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}}.$$

Substituting this expression for  $\alpha_k$  in the formula for  $\gamma_k$  yields

$$\gamma_k = \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)})(\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)})}.$$

Note that in this case  $\gamma_k > 0$  for all  $k$ . Furthermore, by Lemma 8.2, we have  $\gamma_k \geq (\lambda_{\min}(\mathbf{Q})/\lambda_{\max}(\mathbf{Q})) > 0$ . Therefore, we have  $\sum_{k=0}^{\infty} \gamma_k = \infty$ , and hence by Theorem 8.1 we conclude that  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ . ■

Consider now a gradient method with fixed step size; that is,  $\alpha_k = \alpha \in \mathbb{R}$  for all  $k$ . The resulting algorithm is of the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)}.$$

We refer to the algorithm above as a *fixed-step-size* gradient algorithm. The algorithm is of practical interest because of its simplicity. In particular, the algorithm does not require a line search at each step to determine  $\alpha_k$ , because the same step size  $\alpha$  is used at each step. Clearly, the convergence of the algorithm depends on the choice of  $\alpha$ , and we would not expect the algorithm to work for arbitrary  $\alpha$ . The following theorem gives a necessary and sufficient condition on  $\alpha$  for convergence of the algorithm.

**Theorem 8.3** *For the fixed-step-size gradient algorithm,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for any  $\mathbf{x}^{(0)}$  if and only if*

$$0 < \alpha < \frac{2}{\lambda_{\max}(\mathbf{Q})}.$$

□

*Proof.*  $\Leftarrow$ : By Rayleigh's inequality we have

$$\lambda_{\min}(\mathbf{Q})\mathbf{g}^{(k)\top}\mathbf{g}^{(k)} \leq \mathbf{g}^{(k)\top}\mathbf{Q}\mathbf{g}^{(k)} \leq \lambda_{\max}(\mathbf{Q})\mathbf{g}^{(k)\top}\mathbf{g}^{(k)}$$

and

$$\mathbf{g}^{(k)\top}\mathbf{Q}^{-1}\mathbf{g}^{(k)} \leq \frac{1}{\lambda_{\min}(\mathbf{Q})}\mathbf{g}^{(k)\top}\mathbf{g}^{(k)}.$$

Therefore, substituting the above into the formula for  $\gamma_k$ , we get

$$\gamma_k \geq \alpha (\lambda_{\min}(\mathbf{Q}))^2 \left( \frac{2}{\lambda_{\max}(\mathbf{Q})} - \alpha \right) > 0.$$

Therefore,  $\gamma_k > 0$  for all  $k$ , and  $\sum_{k=0}^{\infty} \gamma_k = \infty$ . Hence, by Theorem 8.1 we conclude that  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ .

$\Rightarrow$ : We use contraposition. Suppose that either  $\alpha \leq 0$  or  $\alpha \geq 2/\lambda_{\max}(\mathbf{Q})$ . Let  $\mathbf{x}^{(0)}$  be chosen such that  $\mathbf{x}^{(0)} - \mathbf{x}^*$  is an eigenvector of  $\mathbf{Q}$  corresponding to the eigenvalue  $\lambda_{\max}(\mathbf{Q})$ . Because

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha(\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}) = \mathbf{x}^{(k)} - \alpha(\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{Q}\mathbf{x}^*),$$

we obtain

$$\begin{aligned} \mathbf{x}^{(k+1)} - \mathbf{x}^* &= \mathbf{x}^{(k)} - \mathbf{x}^* - \alpha(\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{Q}\mathbf{x}^*) \\ &= (\mathbf{I}_n - \alpha\mathbf{Q})(\mathbf{x}^{(k)} - \mathbf{x}^*) \\ &= (\mathbf{I}_n - \alpha\mathbf{Q})^{k+1}(\mathbf{x}^{(0)} - \mathbf{x}^*) \\ &= (1 - \alpha\lambda_{\max}(\mathbf{Q}))^{k+1}(\mathbf{x}^{(0)} - \mathbf{x}^*), \end{aligned}$$

where in the last line we used the property that  $\mathbf{x}^{(0)} - \mathbf{x}^*$  is an eigenvector of  $\mathbf{Q}$ . Taking norms on both sides, we get

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = |1 - \alpha\lambda_{\max}(\mathbf{Q})|^{k+1} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|.$$

Because  $\alpha \leq 0$  or  $\alpha \geq 2/\lambda_{\max}(\mathbf{Q})$ ,

$$|1 - \alpha\lambda_{\max}(\mathbf{Q})| \geq 1.$$

Hence,  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|$  cannot converge to 0, and thus the sequence  $\{\mathbf{x}^{(k)}\}$  does not converge to  $\mathbf{x}^*$ . ■

**Example 8.4** Let the function  $f$  be given by

$$f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 4 & 2\sqrt{2} \\ 0 & 5 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 6 \end{bmatrix} + 24.$$

We wish to find the minimizer of  $f$  using a fixed-step-size gradient algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}),$$

where  $\alpha \in \mathbb{R}$  is a fixed step size.

To apply Theorem 8.3, we first symmetrize the matrix in the quadratic term of  $f$  to get

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 8 & 2\sqrt{2} \\ 2\sqrt{2} & 10 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ 6 \end{bmatrix} + 24.$$

The eigenvalues of the matrix in the quadratic term are 6 and 12. Hence, using Theorem 8.3, the algorithm converges to the minimizer for all  $\mathbf{x}^{(0)}$  if and only if  $\alpha$  lies in the range  $0 < \alpha < 2/12$ . ■

## Convergence Rate

We now turn our attention to the issue of convergence rates of gradient algorithms. In particular, we focus on the steepest descent algorithm. We first present the following theorem.

**Theorem 8.4** *In the method of steepest descent applied to the quadratic function, at every step  $k$  we have*

$$V(\mathbf{x}^{(k+1)}) \leq \frac{\lambda_{\max}(\mathbf{Q}) - \lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})} V(\mathbf{x}^{(k)}).$$

□

*Proof.* In the proof of Theorem 8.2, we showed that  $\gamma_k \geq \lambda_{\min}(\mathbf{Q})/\lambda_{\max}(\mathbf{Q})$ . Therefore,

$$\frac{V(\mathbf{x}^{(k)}) - V(\mathbf{x}^{(k+1)})}{V(\mathbf{x}^{(k)})} = \gamma_k \geq \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})},$$

and the result follows. ■

Theorem 8.4 is relevant to our consideration of the convergence rate of the steepest descent algorithm as follows. Let

$$r = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})} = \|\mathbf{Q}\| \|\mathbf{Q}^{-1}\|,$$

called the *condition number* of  $\mathbf{Q}$ . Then, it follows from Theorem 8.4 that

$$V(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{1}{r}\right) V(\mathbf{x}^{(k)}).$$

The term  $(1 - 1/r)$  plays an important role in the convergence of  $\{V(\mathbf{x}^{(k)})\}$  to 0 (and hence of  $\{\mathbf{x}^{(k)}\}$  to  $\mathbf{x}^*$ ). We refer to  $(1 - 1/r)$  as the *convergence ratio*. Specifically, we see that the smaller the value of  $(1 - 1/r)$ , the smaller  $V(\mathbf{x}^{(k+1)})$  will be relative to  $V(\mathbf{x}^{(k)})$ , and hence the “faster”  $V(\mathbf{x}^{(k)})$  converges to 0, as indicated by the inequality above. The convergence ratio  $(1 - 1/r)$  decreases as  $r$  decreases. If  $r = 1$ , then  $\lambda_{\max}(\mathbf{Q}) = \lambda_{\min}(\mathbf{Q})$ , corresponding to circular contours of  $f$  (see Figure 8.6). In this case the algorithm converges in a single step to the minimizer. As  $r$  increases, the speed of convergence of  $\{V(\mathbf{x}^{(k)})\}$  (and hence of  $\{\mathbf{x}^{(k)}\}$ ) decreases. The increase in  $r$  reflects that fact that the contours of  $f$  are more eccentric (see, e.g., Figure 8.7). We refer the reader to [88, pp. 238, 239] for an alternative approach to the analysis above.

To investigate the convergence properties of  $\{\mathbf{x}^{(k)}\}$  further, we need the following definition.

**Definition 8.1** Given a sequence  $\{\mathbf{x}^{(k)}\}$  that converges to  $\mathbf{x}^*$ , that is,  $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$ , we say that the *order of convergence* is  $p$ , where  $p \in \mathbb{R}$ , if

$$0 < \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} < \infty.$$

If for all  $p > 0$ ,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} = 0,$$

then we say that the order of convergence is  $\infty$ . ■

Note that in the definition above,  $0/0$  should be understood to be 0.

The order of convergence of a sequence is a measure of its rate of convergence; the higher the order, the faster the rate of convergence. The order of convergence is sometimes also called the *rate of convergence* (see, e.g., [96]). If  $p = 1$  (first-order convergence) and  $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| / \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 1$ , we say that the convergence is *sublinear*. If  $p = 1$  and  $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| / \|\mathbf{x}^{(k)} - \mathbf{x}^*\| < 1$ , we say that the convergence is *linear*. If  $p > 1$ , we say that the convergence is *superlinear*. If  $p = 2$  (second-order convergence), we say that the convergence is *quadratic*.

**Example 8.5** 1. Suppose that  $x^{(k)} = 1/k$  and thus  $x^{(k)} \rightarrow 0$ . Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{1/(k+1)}{1/k^p} = \frac{k^p}{k+1}.$$

If  $p < 1$ , the sequence above converges to 0, whereas if  $p > 1$ , it grows to  $\infty$ . If  $p = 1$ , the sequence converges to 1. Hence, the order of convergence is 1 (i.e., we have linear convergence).

2. Suppose that  $x^{(k)} = \gamma^k$ , where  $0 < \gamma < 1$ , and thus  $x^{(k)} \rightarrow 0$ . Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{\gamma^{k+1}}{(\gamma^k)^p} = \gamma^{k+1-kp} = \gamma^{k(1-p)+1}.$$

If  $p < 1$ , the sequence above converges to 0, whereas if  $p > 1$ , it grows to  $\infty$ . If  $p = 1$ , the sequence converges to  $\gamma$  (in fact, remains constant at  $\gamma$ ). Hence, the order of convergence is 1.

3. Suppose that  $x^{(k)} = \gamma^{(q^k)}$ , where  $q > 1$  and  $0 < \gamma < 1$ , and thus  $x^{(k)} \rightarrow 0$ . Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{\gamma^{(q^{k+1})}}{(\gamma^{(q^k)})^p} = \gamma^{(q^{k+1}-pq^k)} = \gamma^{(q-p)q^k}.$$

If  $p < q$ , the sequence above converges to 0, whereas if  $p > q$ , it grows to  $\infty$ . If  $p = q$ , the sequence converges to 1 (in fact, remains constant at 1). Hence, the order of convergence is  $q$ .

4. Suppose that  $x^{(k)} = 1$  for all  $k$ , and thus  $x^{(k)} \rightarrow 1$  trivially. Then,

$$\frac{|x^{(k+1)} - 1|}{|x^{(k)} - 1|^p} = \frac{0}{0^p} = 0$$

for all  $p$ . Hence, the order of convergence is  $\infty$ . ■

The order of convergence can be interpreted using the notion of the order symbol  $O$ , as follows. Recall that  $a = O(h)$  ("big-oh of  $h$ ") if there exists a constant  $c$  such that  $|a| \leq c|h|$  for sufficiently small  $h$ . Then, the order of convergence is *at least*  $p$  if

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p)$$

(see Theorem 8.5 below). For example, the order of convergence is at least 2 if

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2)$$

(this fact is used in the analysis of Newton's algorithm in Chapter 9).

**Theorem 8.5** Let  $\{\mathbf{x}^{(k)}\}$  be a sequence that converges to  $\mathbf{x}^*$ . If

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p),$$

then the order of convergence (if it exists) is at least  $p$ .  $\square$

*Proof.* Let  $s$  be the order of convergence of  $\{\mathbf{x}^{(k)}\}$ . Suppose that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p).$$

Then, there exists  $c$  such that for sufficiently large  $k$ ,

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} \leq c.$$

Hence,

$$\begin{aligned} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} &= \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{p-s} \\ &\leq c \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{p-s}. \end{aligned}$$

Taking limits yields

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} \leq c \lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{p-s}.$$

Because by definition  $s$  is the order of convergence,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^s} > 0.$$

Combining the two inequalities above, we get

$$c \lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^{p-s} > 0.$$

Therefore, because  $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$ , we conclude that  $s \geq p$ ; that is, the order of convergence is at least  $p$ .  $\blacksquare$

By an argument similar to the above, we can show that if

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = o(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p),$$

then the order of convergence (if it exists) strictly exceeds  $p$ .

**Example 8.6** Suppose that we are given a scalar sequence  $\{x^{(k)}\}$  that converges with order of convergence  $p$  and satisfies

$$\lim_{k \rightarrow \infty} \frac{|x^{(k+1)} - 2|}{|x^{(k)} - 2|^3} = 0.$$

The limit of  $\{x^{(k)}\}$  must be 2, because it is clear from the equation that  $|x^{(k+1)} - 2| \rightarrow 0$ . Also, we see that  $|x^{(k+1)} - 2| = o(|x^{(k)} - 2|^3)$ . Hence, we conclude that  $p > 3$ . ■

It turns out that the order of convergence of any convergent sequence cannot be less than 1 (see Exercise 8.3). In the following, we provide an example where the order of convergence of a fixed-step-size gradient algorithm exceeds 1.

**Example 8.7** Consider the problem of finding a minimizer of the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  given by

$$f(x) = x^2 - \frac{x^3}{3}.$$

Suppose that we use the algorithm  $x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)})$  with step size  $\alpha = 1/2$  and initial condition  $x^{(0)} = 1$ . (The notation  $f'$  represents the derivative of  $f$ .)

We first show that the algorithm converges to a local minimizer of  $f$ . Indeed, we have  $f'(x) = 2x - x^2$ . The fixed-step-size gradient algorithm with step size  $\alpha = 1/2$  is therefore given by

$$x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)}) = \frac{1}{2}(x^{(k)})^2.$$

With  $x^{(0)} = 1$ , we can derive the expression  $x^{(k)} = (1/2)^{2^k-1}$ . Hence, the algorithm converges to 0, a strict local minimizer of  $f$ .

Next, we find the order of convergence. Note that

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^2} = \frac{1}{2}.$$

Therefore, the order of convergence is 2. ■

Finally, we show that the steepest descent algorithm has an order of convergence of 1 in the *worst case*; that is, there are cases for which the order of convergence of the steepest descent algorithm is equal to 1. To proceed, we will need the following simple lemma.

**Lemma 8.3** *In the steepest descent algorithm, if  $\mathbf{g}^{(k)} \neq \mathbf{0}$  for all  $k$ , then  $\gamma_k = 1$  if and only if  $\mathbf{g}^{(k)}$  is an eigenvector of  $\mathbf{Q}$ .* □

*Proof.* Suppose that  $\mathbf{g}^{(k)} \neq \mathbf{0}$  for all  $k$ . Recall that for the steepest descent algorithm,

$$\gamma_k = \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)})(\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)})}.$$

Sufficiency is easy to show by verification. To show necessity, suppose that  $\gamma_k = 1$ . Then,  $V(\mathbf{x}^{(k+1)}) = 0$ , which implies that  $\mathbf{x}^{(k+1)} = \mathbf{x}^*$ . Therefore,

$$\mathbf{x}^* = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}.$$

Premultiplying by  $\mathbf{Q}$  and subtracting  $\mathbf{b}$  from both sides yields

$$\mathbf{0} = \mathbf{g}^{(k)} - \alpha_k \mathbf{Q} \mathbf{g}^{(k)},$$

which can be rewritten as

$$\mathbf{Q} \mathbf{g}^{(k)} = \frac{1}{\alpha_k} \mathbf{g}^{(k)}.$$

Hence,  $\mathbf{g}^{(k)}$  is an eigenvector of  $\mathbf{Q}$ . ■

By the lemma, if  $\mathbf{g}^{(k)}$  is not an eigenvector of  $\mathbf{Q}$ , then  $\gamma_k < 1$  (recall that  $\gamma_k$  cannot exceed 1). We use this fact in the proof of the following result on the worst-case order of convergence of the steepest descent algorithm.

**Theorem 8.6** *Let  $\{\mathbf{x}^{(k)}\}$  be a convergent sequence of iterates of the steepest descent algorithm applied to a function  $f$ . Then, the order of convergence of  $\{\mathbf{x}^{(k)}\}$  is 1 in the worst case; that is, there exist a function  $f$  and an initial condition  $\mathbf{x}^{(0)}$  such that the order of convergence of  $\{\mathbf{x}^{(k)}\}$  is equal to 1.  $\square$*

*Proof.* Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a quadratic function with Hessian  $\mathbf{Q}$ . Assume that the maximum and minimum eigenvalues of  $\mathbf{Q}$  satisfy  $\lambda_{\max}(\mathbf{Q}) > \lambda_{\min}(\mathbf{Q})$ . To show that the order of convergence of  $\{\mathbf{x}^{(k)}\}$  is 1, it suffices to show that there exists  $\mathbf{x}^{(0)}$  such that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \geq c \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$$

for some  $c > 0$  (see Exercise 8.2). Indeed, by Rayleigh's inequality,

$$\begin{aligned} V(\mathbf{x}^{(k+1)}) &= \frac{1}{2} (\mathbf{x}^{(k+1)} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x}^{(k+1)} - \mathbf{x}^*) \\ &\leq \frac{\lambda_{\max}(\mathbf{Q})}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2. \end{aligned}$$

Similarly,

$$V(\mathbf{x}^{(k)}) \geq \frac{\lambda_{\min}(\mathbf{Q})}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2.$$

Combining the inequalities above with Lemma 8.1, we obtain

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \geq \sqrt{(1 - \gamma_k) \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})}} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|.$$

Therefore, it suffices to choose  $\mathbf{x}^{(0)}$  such that  $\gamma_k \leq d$  for some  $d < 1$ .

Recall that for the steepest descent algorithm, assuming that  $\mathbf{g}^{(k)} \neq \mathbf{0}$  for all  $k$ ,  $\gamma_k$  depends on  $\mathbf{g}^{(k)}$  according to

$$\gamma_k = \frac{(\mathbf{g}^{(k)\top} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)})(\mathbf{g}^{(k)\top} \mathbf{Q}^{-1} \mathbf{g}^{(k)})}.$$

First consider the case where  $n = 2$ . Suppose that  $\mathbf{x}^{(0)} \neq \mathbf{x}^*$  is chosen such that  $\mathbf{x}^{(0)} - \mathbf{x}^*$  is not an eigenvector of  $\mathbf{Q}$ . Then,  $\mathbf{g}^{(0)} = \mathbf{Q}(\mathbf{x}^{(0)} - \mathbf{x}^*) \neq \mathbf{0}$  is also not an eigenvector of  $\mathbf{Q}$ . By Proposition 8.1,  $\mathbf{g}^{(k)} = (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})/\alpha_k$  is not an eigenvector of  $\mathbf{Q}$  for any  $k$  [because any two eigenvectors corresponding to  $\lambda_{\max}(\mathbf{Q})$  and  $\lambda_{\min}(\mathbf{Q})$  are mutually orthogonal]. Also,  $\mathbf{g}^{(k)}$  lies in one of two mutually orthogonal directions. Therefore, by Lemma 8.3, for each  $k$ , the value of  $\gamma_k$  is one of two numbers, both of which are strictly less than 1. This proves the  $n = 2$  case.

For the general  $n$  case, let  $\mathbf{v}_1$  and  $\mathbf{v}_2$  be mutually orthogonal eigenvectors corresponding to  $\lambda_{\max}(\mathbf{Q})$  and  $\lambda_{\min}(\mathbf{Q})$ . Choose  $\mathbf{x}^{(0)}$  such that  $\mathbf{x}^{(0)} - \mathbf{x}^* \neq \mathbf{0}$  lies in the span of  $\mathbf{v}_1$  and  $\mathbf{v}_2$  but is not equal to either. Note that  $\mathbf{g}^{(0)} = \mathbf{Q}(\mathbf{x}^{(0)} - \mathbf{x}^*)$  also lies in the span of  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , but is not equal to either. By manipulating  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$  as before, we can write  $\mathbf{g}^{(k+1)} = (\mathbf{I} - \alpha_k \mathbf{Q})\mathbf{g}^{(k)}$ . Any eigenvector of  $\mathbf{Q}$  is also an eigenvector of  $\mathbf{I} - \alpha_k \mathbf{Q}$ . Therefore,  $\mathbf{g}^{(k)}$  lies in the span of  $\mathbf{v}_1$  and  $\mathbf{v}_2$  for all  $k$ ; that is, the sequence  $\{\mathbf{g}^{(k)}\}$  is confined within the two-dimensional subspace spanned by  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . We can now proceed as in the  $n = 2$  case. ■

In the next chapter we discuss Newton's method, which has order of convergence at least 2 if the initial guess is near the solution.

## EXERCISES

**8.1** Perform two iterations leading to the minimization of

$$f(x_1, x_2) = x_1 + \frac{1}{2}x_2 + \frac{1}{2}x_1^2 + x_2^2 + 3$$

using the steepest descent method with the starting point  $\mathbf{x}^{(0)} = \mathbf{0}$ . Also determine an optimal solution analytically.

**8.2** Let  $\{\mathbf{x}^{(k)}\}$  be a sequence that converges to  $\mathbf{x}^*$ . Show that if there exists  $c > 0$  such that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \geq c\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p$$

for sufficiently large  $k$ , then the order of convergence (if it exists) is at most  $p$ .

**8.3** Let  $\{\mathbf{x}^{(k)}\}$  be a sequence that converges to  $\mathbf{x}^*$ . Show that there does not exist  $p < 1$  such that

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} > 0.$$

**8.4** Consider the sequence  $\{x^{(k)}\}$  given by  $x^{(k)} = 2^{-2^{k^2}}$ .

- a. Write down the value of the limit of  $\{x^{(k)}\}$ .
- b. Find the order of convergence of  $\{x^{(k)}\}$ .

**8.5** Consider the two sequences  $\{x^{(k)}\}$  and  $\{y^{(k)}\}$  defined iteratively as follows:

$$\begin{aligned}x^{(k+1)} &= ax^{(k)}, \\y^{(k+1)} &= (y^{(k)})^b,\end{aligned}$$

where  $a \in \mathbb{R}$ ,  $b \in \mathbb{R}$ ,  $0 < a < 1$ ,  $b > 1$ ,  $x^{(0)} \neq 0$ ,  $y^{(0)} \neq 0$ , and  $|y^{(0)}| < 1$ .

- a. Derive a formula for  $x^{(k)}$  in terms of  $x^{(0)}$  and  $a$ . Use this to deduce that  $x^{(k)} \rightarrow 0$ .
- b. Derive a formula for  $y^{(k)}$  in terms of  $y^{(0)}$  and  $b$ . Use this to deduce that  $y^{(k)} \rightarrow 0$ .
- c. Find the order of convergence of  $\{x^{(k)}\}$  and the order of convergence of  $\{y^{(k)}\}$ .
- d. Calculate the smallest number of iterations  $k$  such that  $|x^{(k)}| \leq c|x^{(0)}|$ , where  $0 < c < 1$ .  
*Hint:* The answer is in terms of  $a$  and  $c$ . You may use the notation  $[z]$  to represent the smallest integer not smaller than  $z$ .
- e. Calculate the smallest number of iterations  $k$  such that  $|y^{(k)}| \leq c|y^{(0)}|$ , where  $0 < c < 1$ .
- f. Compare the answer of part e with that of part d, focusing on the case where  $c$  is very small.

**8.6** Suppose that we use the golden section algorithm to find the minimizer of a function. Let  $u_k$  be the uncertainty range at the  $k$ th iteration. Find the order of convergence of  $\{u_k\}$ .

**8.7** Suppose that we wish to minimize a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  that has a derivative  $f'$ . A simple line search method, called *derivative descent search* (DDS), is described as follows: given that we are at a point  $x^{(k)}$ , we move in the direction of the negative derivative with step size  $\alpha$ ; that is,  $x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)})$ , where  $\alpha > 0$  is a constant.

In the following parts, assume that  $f$  is quadratic:  $f(x) = \frac{1}{2}ax^2 - bx + c$  (where  $a$ ,  $b$ , and  $c$  are constants, and  $a > 0$ ).

- a. Write down the value of  $x^*$  (in terms of  $a$ ,  $b$ , and  $c$ ) that minimizes  $f$ .

- b. Write down the recursive equation for the DDS algorithm explicitly for this quadratic  $f$ .
- c. Assuming that the DDS algorithm converges, show that it converges to the optimal value  $x^*$  (found in part a).
- d. Find the order of convergence of the algorithm, assuming that it does converge.
- e. Find the range of values of  $\alpha$  for which the algorithm converges (for this particular  $f$ ) for all starting points  $x^{(0)}$ .

**8.8** Consider the function

$$f(\mathbf{x}) = 3(x_1^2 + x_2^2) + 4x_1x_2 + 5x_1 + 6x_2 + 7,$$

where  $\mathbf{x} = [x_1, x_2]^\top \in \mathbb{R}^2$ . Suppose that we use a fixed-step-size gradient algorithm to find the minimizer of  $f$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}).$$

Find the largest range of values of  $\alpha$  for which the algorithm is globally convergent.

**8.9** This exercise explores a zero-finding algorithm.

Suppose that we wish to solve the equation  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ , where

$$\mathbf{h}(\mathbf{x}) = \begin{bmatrix} 4 + 3x_1 + 2x_2 \\ 1 + 2x_1 + 3x_2 \end{bmatrix}.$$

Consider using an algorithm of the form  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{h}(\mathbf{x}^{(k)})$ , where  $\alpha$  is scalar constant that does not depend on  $k$ .

- a. Find the solution of  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ .
- b. Find the largest range of values of  $\alpha$  such that the algorithm is globally convergent to the solution of  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ .
- c. Assuming that  $\alpha$  is outside the range of values in part b, give an example of an initial condition  $\mathbf{x}^{(0)}$  of the form  $[x_1, 0]^\top$  such that the algorithm is guaranteed not to satisfy the descent property.

**8.10** Consider the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$f(\mathbf{x}) = \frac{3}{2}(x_1^2 + x_2^2) + (1+a)x_1x_2 - (x_1 + x_2) + b,$$

where  $a$  and  $b$  are some unknown real-valued parameters.

- a. Write the function  $f$  in the usual multivariable quadratic form.
- b. Find the largest set of values of  $a$  and  $b$  such that the unique global minimizer of  $f$  exists, and write down the minimizer (in terms of the parameters  $a$  and  $b$ ).
- c. Consider the following algorithm:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{2}{5} \nabla f(\mathbf{x}^{(k)}).$$

Find the largest set of values of  $a$  and  $b$  for which this algorithm converges to the global minimizer of  $f$  for any initial point  $\mathbf{x}^{(0)}$ .

**8.11** Consider the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  given by  $f(x) = \frac{1}{2}(x - c)^2$ ,  $c \in \mathbb{R}$ . We are interested in computing the minimizer of  $f$  using the iterative algorithm

$$x^{(k+1)} = x^{(k)} - \alpha_k f'(x^{(k)}),$$

where  $f'$  is the derivative of  $f$  and  $\alpha_k$  is a step size satisfying  $0 < \alpha_k < 1$ .

- a. Derive a formula relating  $f(x^{(k+1)})$  with  $f(x^{(k)})$ , involving  $\alpha_k$ .
- b. Show that the algorithm is globally convergent if and only if

$$\sum_{k=0}^{\infty} \alpha_k = \infty.$$

*Hint:* Use part a and the fact that for any sequence  $\{\alpha_k\} \subset (0, 1)$ , we have

$$\prod_{k=0}^{\infty} (1 - \alpha_k) = 0 \Leftrightarrow \sum_{k=0}^{\infty} \alpha_k = \infty.$$

**8.12** Consider the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  given by  $f(x) = x^3 - x$ . Suppose that we use a fixed-step-size algorithm  $x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)})$  to find a local minimizer of  $f$ . Find the largest range of values of  $\alpha$  such that the algorithm is locally convergent (i.e., for all  $x_0$  sufficiently close to a local minimizer  $x^*$ , we have  $x^{(k)} \rightarrow x^*$ ).

**8.13** Consider the function  $f$  given by  $f(x) = (x - 1)^2$ ,  $x \in \mathbb{R}$ . We are interested in computing the minimizer of  $f$  using the iterative algorithm  $x^{(k+1)} = x^{(k)} - \alpha 2^{-k} f'(x^{(k)})$ , where  $f'$  is the derivative of  $f$  and  $0 < \alpha < 1$ . Does the algorithm have the descent property? Is the algorithm globally convergent?

**8.14** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^3$ , with first derivative  $f'$ , second derivative  $f''$ , and unique minimizer  $x^*$ . Consider a fixed-step-size gradient algorithm

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

Suppose that  $f''(x^*) \neq 0$  and  $\alpha = 1/f''(x^*)$ . Assuming that the algorithm converges to  $x^*$ , show that the order of convergence is at least 2.

**8.15** Consider the problem of minimizing  $f(x) = \|\mathbf{ax} - \mathbf{b}\|^2$ , where  $\mathbf{a}$  and  $\mathbf{b}$  are vectors in  $\mathbb{R}^n$ , and  $\mathbf{a} \neq \mathbf{0}$ .

- a. Derive an expression (in terms of  $\mathbf{a}$  and  $\mathbf{b}$ ) for the solution to this problem.
- b. To solve the problem, suppose that we use an iterative algorithm of the form

$$x^{(k+1)} = x^{(k)} - \alpha f'(x^{(k)}),$$

where  $f'$  is the derivative of  $f$ . Find the largest range of values of  $\alpha$  (in terms of  $\mathbf{a}$  and  $\mathbf{b}$ ) for which the algorithm converges to the solution for all starting points  $x^{(0)}$ .

**8.16** Consider the optimization problem

$$\text{minimize } \|\mathbf{Ax} - \mathbf{b}\|^2,$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ , and  $\mathbf{b} \in \mathbb{R}^m$ .

- a. Show that the objective function for this problem is a quadratic function, and write down the gradient and Hessian of this quadratic.
- b. Write down the fixed-step-size gradient algorithm for solving this optimization problem.
- c. Suppose that

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}.$$

Find the largest range of values for  $\alpha$  such that the algorithm in part b converges to the solution of the problem.

**8.17** Consider a function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  given by  $\mathbf{f}(\mathbf{x}) = \mathbf{Ax} + \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ . Suppose that  $\mathbf{A}$  is invertible and  $\mathbf{x}^*$  is the zero of  $\mathbf{f}$  [i.e.,  $\mathbf{f}(\mathbf{x}^*) = \mathbf{0}$ ]. We wish to compute  $\mathbf{x}^*$  using the iterative algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{f}(\mathbf{x}^{(k)}),$$

where  $\alpha \in \mathbb{R}$ ,  $\alpha > 0$ . We say that the algorithm is *globally monotone* if for any  $\mathbf{x}^{(0)}$ ,  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$  for all  $k$ .

- a. Assume that all the eigenvalues of  $\mathbf{A}$  are real. Show that a necessary condition for the algorithm above to be *globally monotone* is that all the eigenvalues of  $\mathbf{A}$  are nonnegative.

*Hint:* Use contraposition.

- b. Suppose that

$$\mathbf{A} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ -1 \end{bmatrix}.$$

Find the largest range of values of  $\alpha$  for which the algorithm is *globally convergent* (i.e.,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$  for all  $\mathbf{x}^{(0)}$ ).

- 8.18** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , where  $\mathbf{b} \in \mathbb{R}^n$  and  $\mathbf{Q}$  is a real symmetric positive definite  $n \times n$  matrix. Suppose that we apply the steepest descent method to this function, with  $\mathbf{x}^{(0)} \neq \mathbf{Q}^{-1}\mathbf{b}$ . Show that the method converges in one step, that is,  $\mathbf{x}^{(1)} = \mathbf{Q}^{-1}\mathbf{b}$ , if and only if  $\mathbf{x}^{(0)}$  is chosen such that  $\mathbf{g}^{(0)} = \mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}$  is an eigenvector of  $\mathbf{Q}$ .

- 8.19** Suppose that we apply the steepest descent algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$  to a quadratic function  $f$  with Hessian  $\mathbf{Q} > 0$ . Let  $\lambda_{\max}$  and  $\lambda_{\min}$  be the largest and smallest eigenvalue of  $\mathbf{Q}$ , respectively. Which of the following two inequalities are possibly true? (When we say here that an inequality is “possibly” true, we mean that there exists a choice of  $f$  and  $\mathbf{x}^{(0)}$  such that the inequality holds.)

- a.  $\alpha_0 \geq 2/\lambda_{\max}$ .
- b.  $\alpha_0 > 1/\lambda_{\min}$ .

- 8.20** Suppose that we apply a fixed-step-size gradient algorithm to minimize

$$f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 3/2 & 2 \\ 0 & 3/2 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 3 \\ -1 \end{bmatrix} - 22.$$

- a. Find the range of values of the step size for which the algorithm converges to the minimizer.
- b. Suppose that we use a step size of 1000 (which is too large). Find an initial condition that will cause the algorithm to diverge (not converge).

- 8.21** Consider a fixed-step-size gradient algorithm applied to each of the functions  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  in parts a and b below. In each case, find the largest range of values of the step size  $\alpha$  for which the algorithm is globally convergent.

- a.  $f(\mathbf{x}) = 1 + 2x_1 + 3(x_1^2 + x_2^2) + 4x_1x_2$ .

b.  $f(\mathbf{x}) = \mathbf{x}^\top \begin{bmatrix} 3 & 3 \\ 1 & 3 \end{bmatrix} \mathbf{x} + [16, 23]\mathbf{x} + \pi^2.$

**8.22** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , where  $\mathbf{b} \in \mathbb{R}^n$  and  $\mathbf{Q}$  is a real symmetric positive definite  $n \times n$  matrix. Consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \beta \alpha_k \mathbf{g}^{(k)},$$

where  $\mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}$ ,  $\alpha_k = \mathbf{g}^{(k)\top} \mathbf{g}^{(k)} / \mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}$ , and  $\beta \in \mathbb{R}$  is a given constant. (Note that the above reduces to the steepest descent algorithm if  $\beta = 1$ .) Show that  $\{\mathbf{x}^{(k)}\}$  converges to  $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{b}$  for any initial condition  $\mathbf{x}^{(0)}$  if and only if  $0 < \beta < 2$ .

**8.23** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , where  $\mathbf{b} \in \mathbb{R}^n$  and  $\mathbf{Q}$  is a real symmetric positive definite  $n \times n$  matrix. Consider a gradient algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)},$$

where  $\mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}$  is the gradient of  $f$  at  $\mathbf{x}^{(k)}$  and  $\alpha_k$  is some step size. Show that the algorithm has the descent property [i.e.,  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$  whenever  $\mathbf{g}^{(k)} \neq \mathbf{0}$ ] if and only if  $\gamma_k > 0$  for all  $k$ .

**8.24** Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , consider the general iterative algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

where  $\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots$  are given vectors in  $\mathbb{R}^n$  and  $\alpha_k$  is chosen to minimize  $f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ ; that is,

$$\alpha_k = \arg \min f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

Show that for each  $k$ , the vector  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$  is orthogonal to  $\nabla f(\mathbf{x}^{(k+1)})$  (assuming that the gradient exists).

**8.25** Write a simple MATLAB program for implementing the steepest descent algorithm using the secant method for the line search (e.g., the MATLAB function of Exercise 7.11). For the stopping criterion, use the condition  $\|\mathbf{g}^{(k)}\| \leq \varepsilon$ , where  $\varepsilon = 10^{-6}$ . Test your program by comparing the output with the numbers in Example 8.1. Also test your program using an initial condition of  $[-4, 5, 1]^\top$ , and determine the number of iterations required to satisfy the stopping criterion. Evaluate the objective function at the final point to see how close it is to 0.

**8.26** Apply the MATLAB program from Exercise 8.25 to Rosenbrock's function:

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2.$$

Use an initial condition of  $\mathbf{x}^{(0)} = [-2, 2]^\top$ . Terminate the algorithm when the norm of the gradient of  $f$  is less than  $10^{-4}$ .



# CHAPTER 9

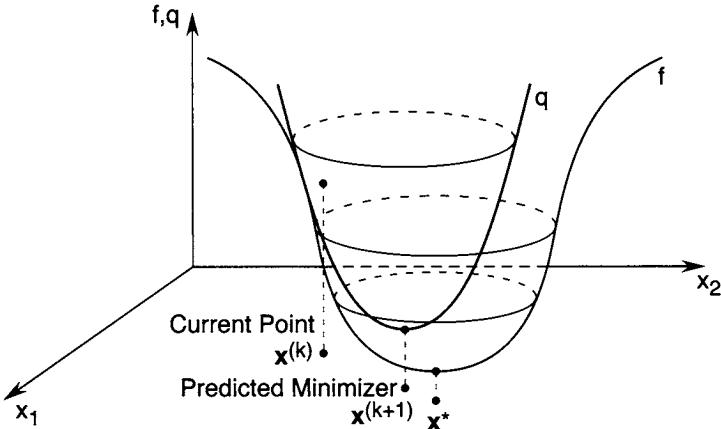
---

## NEWTON'S METHOD

---

### 9.1 Introduction

Recall that the method of steepest descent uses only first derivatives (gradients) in selecting a suitable search direction. This strategy is not always the most effective. If higher derivatives are used, the resulting iterative algorithm may perform better than the steepest descent method. *Newton's method* (sometimes called the *Newton-Raphson method*) uses first and second derivatives and indeed does perform better than the steepest descent method if the initial point is close to the minimizer. The idea behind this method is as follows. Given a starting point, we construct a quadratic approximation to the objective function that matches the first and second derivative values at that point. We then minimize the approximate (quadratic) function instead of the original objective function. We use the minimizer of the approximate function as the starting point in the next step and repeat the procedure iteratively. If the objective function is quadratic, then the approximation is exact, and the method yields the true minimizer in one step. If, on the other hand, the objective function is not quadratic, then the approximation will provide



**Figure 9.1** Quadratic approximation to the objective function using first and second derivatives.

only an estimate of the position of the true minimizer. Figure 9.1 illustrates this idea.

We can obtain a quadratic approximation to the twice continuously differentiable objection function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  using the Taylor series expansion of  $f$  about the current point  $\mathbf{x}^{(k)}$ , neglecting terms of order three and higher. We obtain

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{g}^{(k)} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{F}(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) \triangleq q(\mathbf{x}),$$

where, for simplicity, we use the notation  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ . Applying the FONC to  $q$  yields

$$\mathbf{0} = \nabla q(\mathbf{x}) = \mathbf{g}^{(k)} + \mathbf{F}(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}).$$

If  $\mathbf{F}(\mathbf{x}^{(k)}) > 0$ , then  $q$  achieves a minimum at

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}.$$

This recursive formula represents Newton's method.

**Example 9.1** Use Newton's method to minimize the Powell function:

$$f(x_1, x_2, x_3, x_4) = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 - 2x_3)^4 + 10(x_1 - x_4)^4.$$

Use as the starting point  $\mathbf{x}^{(0)} = [3, -1, 0, 1]^\top$ . Perform three iterations.

Note that  $f(\mathbf{x}^{(0)}) = 215$ . We have

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 2(x_1 + 10x_2) + 40(x_1 - x_4)^3 \\ 20(x_1 + 10x_2) + 4(x_2 - 2x_3)^3 \\ 10(x_3 - x_4) - 8(x_2 - 2x_3)^3 \\ -10(x_3 - x_4) - 40(x_1 - x_4)^3 \end{bmatrix},$$

and  $\mathbf{F}(\mathbf{x})$  is given by

$$\begin{bmatrix} 2 + 120(x_1 - x_4)^2 & 20 & 0 & -120(x_1 - x_4)^2 \\ 20 & 200 + 12(x_2 - 2x_3)^2 & -24(x_2 - 2x_3)^2 & 0 \\ 0 & -24(x_2 - 2x_3)^2 & 10 + 48(x_2 - 2x_3)^2 & -10 \\ -120(x_1 - x_4)^2 & 0 & -10 & 10 + 120(x_1 - x_4)^2 \end{bmatrix}.$$

*Iteration 1*

$$\begin{aligned} \mathbf{g}^{(0)} &= [306, -144, -2, -310]^\top, \\ \mathbf{F}(\mathbf{x}^{(0)}) &= \begin{bmatrix} 482 & 20 & 0 & -480 \\ 20 & 212 & -24 & 0 \\ 0 & -24 & 58 & -10 \\ -480 & 0 & -10 & 490 \end{bmatrix}, \\ \mathbf{F}(\mathbf{x}^{(0)})^{-1} &= \begin{bmatrix} 0.1126 & -0.0089 & 0.0154 & 0.1106 \\ -0.0089 & 0.0057 & 0.0008 & -0.0087 \\ 0.0154 & 0.0008 & 0.0203 & 0.0155 \\ 0.1106 & -0.0087 & 0.0155 & 0.1107 \end{bmatrix}, \\ \mathbf{F}(\mathbf{x}^{(0)})^{-1}\mathbf{g}^{(0)} &= [1.4127, -0.8413, -0.2540, 0.7460]^\top. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)})^{-1}\mathbf{g}^{(0)} = [1.5873, -0.1587, 0.2540, 0.2540]^\top, \\ f(\mathbf{x}^{(1)}) &= 31.8. \end{aligned}$$

*Iteration 2*

$$\begin{aligned} \mathbf{g}^{(1)} &= [94.81, -1.179, 2.371, -94.81]^\top, \\ \mathbf{F}(\mathbf{x}^{(1)}) &= \begin{bmatrix} 215.3 & 20 & 0 & -213.3 \\ 20 & 205.3 & -10.67 & 0 \\ 0 & -10.67 & 31.34 & -10 \\ -213.3 & 0 & -10 & 223.3 \end{bmatrix}, \\ \mathbf{F}(\mathbf{x}^{(1)})^{-1}\mathbf{g}^{(1)} &= [0.5291, -0.0529, 0.0846, 0.0846]^\top. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbf{x}^{(2)} &= \mathbf{x}^{(1)} - \mathbf{F}(\mathbf{x}^{(1)})^{-1}\mathbf{g}^{(1)} = [1.0582, -0.1058, 0.1694, 0.1694]^\top, \\ f(\mathbf{x}^{(2)}) &= 6.28. \end{aligned}$$

*Iteration 3*

$$\begin{aligned}\mathbf{g}^{(2)} &= [28.09, -0.3475, 0.7031, -28.08]^\top, \\ \mathbf{F}(\mathbf{x}^{(2)}) &= \begin{bmatrix} 96.80 & 20 & 0 & -94.80 \\ 20 & 202.4 & -4.744 & 0 \\ 0 & -4.744 & 19.49 & -10 \\ -94.80 & 0 & -10 & 104.80 \end{bmatrix}, \\ \mathbf{x}^{(3)} &= [0.7037, -0.0704, 0.1121, 0.1111]^\top, \\ f(\mathbf{x}^{(3)}) &= 1.24.\end{aligned}$$

■

Observe that the  $k$ th iteration of Newton's method can be written in two steps as

1. Solve  $\mathbf{F}(\mathbf{x}^{(k)})\mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$  for  $\mathbf{d}^{(k)}$ .
2. Set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$ .

Step 1 requires the solution of an  $n \times n$  system of linear equations. Thus, an efficient method for solving systems of linear equations is essential when using Newton's method.

As in the one-variable case, Newton's method can also be viewed as a technique for iteratively solving the equation

$$\mathbf{g}(\mathbf{x}) = \mathbf{0},$$

where  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . In this case  $\mathbf{F}(\mathbf{x})$  is the Jacobian matrix of  $\mathbf{g}$  at  $\mathbf{x}$ ; that is,  $\mathbf{F}(\mathbf{x})$  is the  $n \times n$  matrix whose  $(i, j)$  entry is  $(\partial g_i / \partial x_j)(\mathbf{x})$ ,  $i, j = 1, 2, \dots, n$ .

## 9.2 Analysis of Newton's Method

As in the one-variable case there is no guarantee that Newton's algorithm heads in the direction of decreasing values of the objective function if  $\mathbf{F}(\mathbf{x}^{(k)})$  is not positive definite (recall Figure 7.7 illustrating Newton's method for functions of one variable when  $f'' < 0$ ). Moreover, even if  $\mathbf{F}(\mathbf{x}^{(k)}) > 0$ , Newton's method may not be a descent method; that is, it is possible that  $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$ . For example, this may occur if our starting point  $\mathbf{x}^{(0)}$  is far away from the solution. See the end of this section for a possible remedy to this problem. Despite these drawbacks, Newton's method has superior convergence properties when the starting point is near the solution, as we shall see in the remainder of this section.

The convergence analysis of Newton's method when  $f$  is a quadratic function is straightforward. In fact, Newton's method reaches the point  $\mathbf{x}^*$  such

that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  in just one step starting from any initial point  $\mathbf{x}^{(0)}$ . To see this, suppose that  $\mathbf{Q} = \mathbf{Q}^\top$  is invertible and

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}.$$

Then,

$$\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b}$$

and

$$\mathbf{F}(\mathbf{x}) = \mathbf{Q}.$$

Hence, given any initial point  $\mathbf{x}^{(0)}$ , by Newton's algorithm

$$\begin{aligned}\mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)})^{-1}\mathbf{g}^{(0)} \\ &= \mathbf{x}^{(0)} - \mathbf{Q}^{-1}[\mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}] \\ &= \mathbf{Q}^{-1}\mathbf{b} \\ &= \mathbf{x}^*.\end{aligned}$$

Therefore, for the quadratic case the order of convergence of Newton's algorithm is  $\infty$  for any initial point  $\mathbf{x}^{(0)}$  (compare this with Exercise 8.18, which deals with the steepest descent algorithm).

To analyze the convergence of Newton's method in the general case, we use results from Section 5.1. Let  $\{\mathbf{x}^{(k)}\}$  be the Newton's method sequence for minimizing a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . We show that  $\{\mathbf{x}^{(k)}\}$  converges to the minimizer  $\mathbf{x}^*$  with order of convergence at least 2.

**Theorem 9.1** Suppose that  $f \in \mathcal{C}^3$  and  $\mathbf{x}^* \in \mathbb{R}^n$  is a point such that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and  $\mathbf{F}(\mathbf{x}^*)$  is invertible. Then, for all  $\mathbf{x}^{(0)}$  sufficiently close to  $\mathbf{x}^*$ , Newton's method is well-defined for all  $k$  and converges to  $\mathbf{x}^*$  with an order of convergence at least 2.  $\square$

*Proof.* The Taylor series expansion of  $\nabla f$  about  $\mathbf{x}^{(0)}$  yields

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)}) = O(\|\mathbf{x} - \mathbf{x}^{(0)}\|^2).$$

Because by assumption  $f \in \mathcal{C}^3$  and  $\mathbf{F}(\mathbf{x}^*)$  is invertible, there exist constants  $\varepsilon > 0$ ,  $c_1 > 0$ , and  $c_2 > 0$  such that if  $\mathbf{x}^{(0)}, \mathbf{x} \in \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \varepsilon\}$ , we have

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)})\| \leq c_1 \|\mathbf{x} - \mathbf{x}^{(0)}\|^2$$

and by Lemma 5.3,  $\mathbf{F}(\mathbf{x})^{-1}$  exists and satisfies

$$\|\mathbf{F}(\mathbf{x})^{-1}\| \leq c_2.$$

The first inequality above holds because the remainder term in the Taylor series expansion contains third derivatives of  $f$  that are continuous and hence bounded on  $\{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \varepsilon\}$ .

Suppose that  $\mathbf{x}^{(0)} \in \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \varepsilon\}$ . Then, substituting  $\mathbf{x} = \mathbf{x}^*$  in the inequality above and using the assumption that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ , we get

$$\|\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\| \leq c_1 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2.$$

Now, subtracting  $\mathbf{x}^*$  from both sides of Newton's algorithm and taking norms yields

$$\begin{aligned}\|\mathbf{x}^{(1)} - \mathbf{x}^*\| &= \|\mathbf{x}^{(0)} - \mathbf{x}^* - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \nabla f(\mathbf{x}^{(0)})\| \\ &= \|\mathbf{F}(\mathbf{x}^{(0)})^{-1} (\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)}))\| \\ &\leq \|\mathbf{F}(\mathbf{x}^{(0)})^{-1}\| \|\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\|.\end{aligned}$$

Applying the inequalities above involving the constants  $c_1$  and  $c_2$  gives

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2.$$

Suppose that  $\mathbf{x}^{(0)}$  is such that

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \frac{\alpha}{c_1 c_2},$$

where  $\alpha \in (0, 1)$ . Then,

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\| \leq \alpha \|\mathbf{x}^{(0)} - \mathbf{x}^*\|.$$

By induction, we obtain

$$\begin{aligned}\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| &\leq c_1 c_2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2, \\ \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| &\leq \alpha \|\mathbf{x}^{(k)} - \mathbf{x}^*\|.\end{aligned}$$

Hence,

$$\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0,$$

and therefore the sequence  $\{\mathbf{x}^{(k)}\}$  converges to  $\mathbf{x}^*$ . The order of convergence is at least 2 because  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$ ; that is,  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2)$ . ■

**Warning:** In the Theorem 9.1, we did not state that  $\mathbf{x}^*$  is a local minimizer. For example, if  $\mathbf{x}^*$  is a local *maximizer*, then provided that  $f \in C^3$  and  $\mathbf{F}(\mathbf{x}^*)$  is invertible, Newton's method would converge to  $\mathbf{x}^*$  if we start close enough to it.

As stated in Theorem 9.1, Newton's method has superior convergence properties if the starting point is near the solution. However, the method is not guaranteed to converge to the solution if we start far away from it (in fact, it may not even be well-defined because the Hessian may be singular). In particular, the method may not be a descent method; that is, it is possible that

$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$ . Fortunately, it is possible to modify the algorithm such that the descent property holds. To see this, we need the following result.

**Theorem 9.2** *Let  $\{\mathbf{x}^{(k)}\}$  be the sequence generated by Newton's method for minimizing a given objective function  $f(\mathbf{x})$ . If the Hessian  $\mathbf{F}(\mathbf{x}^{(k)}) > 0$  and  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , then the search direction*

$$\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1}\mathbf{g}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$$

*from  $\mathbf{x}^{(k)}$  to  $\mathbf{x}^{(k+1)}$  is a descent direction for  $f$  in the sense that there exists an  $\bar{\alpha} > 0$  such that for all  $\alpha \in (0, \bar{\alpha})$ ,*

$$f(\mathbf{x}^{(k)} + \alpha\mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}).$$

□

*Proof.* Let

$$\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha\mathbf{d}^{(k)}).$$

Then, using the chain rule, we obtain

$$\phi'(\alpha) = \nabla f(\mathbf{x}^{(k)} + \alpha\mathbf{d}^{(k)})^\top \mathbf{d}^{(k)}.$$

Hence,

$$\phi'(0) = \nabla f(\mathbf{x}^{(k)})^\top \mathbf{d}^{(k)} = -\mathbf{g}^{(k)\top} \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} < 0,$$

because  $\mathbf{F}(\mathbf{x}^{(k)})^{-1} > 0$  and  $\mathbf{g}^{(k)} \neq \mathbf{0}$ . Thus, there exists an  $\bar{\alpha} > 0$  so that for all  $\alpha \in (0, \bar{\alpha})$ ,  $\phi(\alpha) < \phi(0)$ . This implies that for all  $\alpha \in (0, \bar{\alpha})$ ,

$$f(\mathbf{x}^{(k)} + \alpha\mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}),$$

which completes the proof. ■

Theorem 9.2 motivates the following modification of Newton's method:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)},$$

where

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)});$$

that is, at each iteration, we perform a line search in the direction  $-\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ . By Theorem 9.2 we conclude that the modified Newton's method has the descent property; that is,

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$$

whenever  $\mathbf{g}^{(k)} \neq \mathbf{0}$ .

A drawback of Newton's method is that evaluation of  $\mathbf{F}(\mathbf{x}^{(k)})$  for large  $n$  can be computationally expensive. Furthermore, we have to solve the set of  $n$  linear equations  $\mathbf{F}(\mathbf{x}^{(k)})\mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$ . In Chapters 10 and 11 we discuss methods that alleviate this difficulty.

Another source of potential problems in Newton's method arises from the Hessian matrix not being positive definite. In the next section we describe a simple modification of Newton's method to overcome this problem.

### 9.3 Levenberg-Marquardt Modification

If the Hessian matrix  $\mathbf{F}(\mathbf{x}^{(k)})$  is not positive definite, then the search direction  $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1}\mathbf{g}^{(k)}$  may not point in a descent direction. A simple technique to ensure that the search direction is a descent direction is to introduce the *Levenberg-Marquardt modification* of Newton's algorithm:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)},$$

where  $\mu_k \geq 0$ .

The idea underlying the Levenberg-Marquardt modification is as follows. Consider a symmetric matrix  $\mathbf{F}$ , which may not be positive definite. Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of  $\mathbf{F}$  with corresponding eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$ . The eigenvalues  $\lambda_1, \dots, \lambda_n$  are real, but may not all be positive. Next, consider the matrix  $\mathbf{G} = \mathbf{F} + \mu \mathbf{I}$ , where  $\mu \geq 0$ . Note that the eigenvalues of  $\mathbf{G}$  are  $\lambda_1 + \mu, \dots, \lambda_n + \mu$ . Indeed,

$$\begin{aligned} \mathbf{G}\mathbf{v}_i &= (\mathbf{F} + \mu \mathbf{I})\mathbf{v}_i \\ &= \mathbf{F}\mathbf{v}_i + \mu \mathbf{I}\mathbf{v}_i \\ &= \lambda_i \mathbf{v}_i + \mu \mathbf{v}_i \\ &= (\lambda_i + \mu) \mathbf{v}_i, \end{aligned}$$

which shows that for all  $i = 1, \dots, n$ ,  $\mathbf{v}_i$  is also an eigenvector of  $\mathbf{G}$  with eigenvalue  $\lambda_i + \mu$ . Therefore, if  $\mu$  is sufficiently large, then all the eigenvalues of  $\mathbf{G}$  are positive and  $\mathbf{G}$  is positive definite. Accordingly, if the parameter  $\mu_k$  in the Levenberg-Marquardt modification of Newton's algorithm is sufficiently large, then the search direction  $\mathbf{d}^{(k)} = -(\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)}$  always points in a descent direction (in the sense of Theorem 9.2). In this case if we further introduce a step size  $\alpha_k$  as described in Section 9.2,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k (\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)},$$

then we are guaranteed that the descent property holds.

The Levenberg-Marquardt modification of Newton's algorithm can be made to approach the behavior of the pure Newton's method by letting  $\mu_k \rightarrow 0$ . On the other hand, by letting  $\mu_k \rightarrow \infty$ , the algorithm approaches a pure gradient method with small step size. In practice, we may start with a small value of  $\mu_k$  and increase it slowly until we find that the iteration is descent:  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ .

### 9.4 Newton's Method for Nonlinear Least Squares

We now examine a particular class of optimization problems and the use of Newton's method for solving them. Consider the following problem:

$$\text{minimize } \sum_{i=1}^m (r_i(\mathbf{x}))^2,$$

where  $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i = 1, \dots, m$ , are given functions. This particular problem is called a *nonlinear least-squares problem*. The special case where the  $r_i$  are linear is discussed in Section 12.1.

**Example 9.2** Suppose that we are given  $m$  measurements of a process at  $m$  points in time, as depicted in Figure 9.2 (here,  $m = 21$ ). Let  $t_1, \dots, t_m$  denote the measurement times and  $y_1, \dots, y_m$  the measurement values. Note that  $t_1 = 0$  while  $t_{21} = 10$ . We wish to fit a sinusoid to the measurement data. The equation of the sinusoid is

$$y = A \sin(\omega t + \phi)$$

with appropriate choices of the parameters  $A$ ,  $\omega$ , and  $\phi$ . To formulate the data-fitting problem, we construct the objective function

$$\sum_{i=1}^m (y_i - A \sin(\omega t_i + \phi))^2,$$

representing the sum of the squared errors between the measurement values and the function values at the corresponding points in time. Let  $\mathbf{x} = [A, \omega, \phi]^\top$  represent the vector of decision variables. We therefore obtain a nonlinear least-squares problem with

$$r_i(\mathbf{x}) = y_i - A \sin(\omega t_i + \phi).$$

■

Defining  $\mathbf{r} = [r_1, \dots, r_m]^\top$ , we write the objective function as  $f(\mathbf{x}) = \mathbf{r}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$ . To apply Newton's method, we need to compute the gradient and the Hessian of  $f$ . The  $j$ th component of  $\nabla f(\mathbf{x})$  is

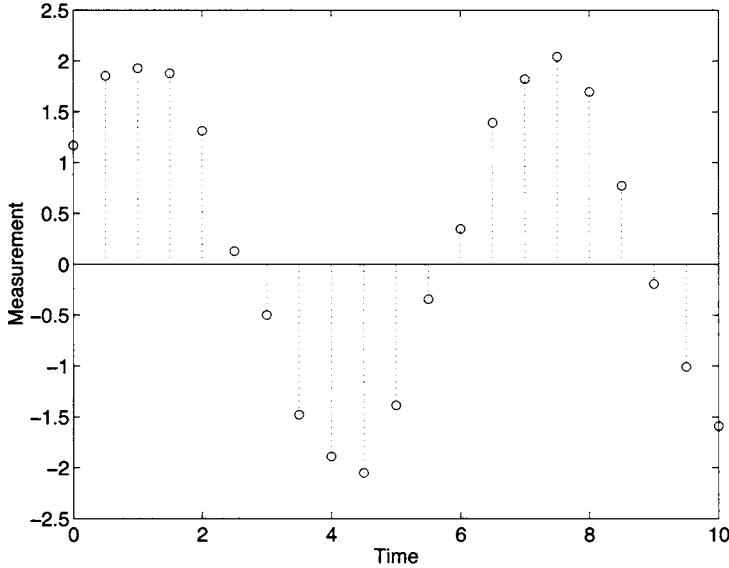
$$(\nabla f(\mathbf{x}))_j = \frac{\partial f}{\partial x_j}(\mathbf{x}) = 2 \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}).$$

Denote the Jacobian matrix of  $\mathbf{r}$  by

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial r_1}{\partial x_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial r_m}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial r_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}.$$

Then, the gradient of  $f$  can be represented as

$$\nabla f(\mathbf{x}) = 2\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$



**Figure 9.2** Measurement data for Example 9.2.

Next, we compute the Hessian matrix of  $f$ . The  $(k, j)$ th component of the Hessian is given by

$$\begin{aligned} \frac{\partial^2 f}{\partial x_k \partial x_j}(\mathbf{x}) &= \frac{\partial}{\partial x_k} \left( \frac{\partial f}{\partial x_j}(\mathbf{x}) \right) \\ &= \frac{\partial}{\partial x_k} \left( 2 \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) \right) \\ &= 2 \sum_{i=1}^m \left( \frac{\partial r_i}{\partial x_k}(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) + r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x}) \right). \end{aligned}$$

Letting  $\mathbf{S}(\mathbf{x})$  be the matrix whose  $(k, j)$ th component is

$$\sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x}),$$

we write the Hessian matrix as

$$\mathbf{F}(\mathbf{x}) = 2(\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mathbf{S}(\mathbf{x})).$$

Therefore, Newton's method applied to the nonlinear least-squares problem is given by

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mathbf{S}(\mathbf{x}))^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$

In some applications, the matrix  $\mathbf{S}(\mathbf{x})$  involving the second derivatives of the function  $\mathbf{r}$  can be ignored because its components are negligibly small. In this case Newton's algorithm reduces to what is commonly called the *Gauss-Newton method*:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}))^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$

Note that the Gauss-Newton method does not require calculation of the second derivatives of  $\mathbf{r}$ .

**Example 9.3** Recall the data-fitting problem in Example 9.2, with

$$r_i(\mathbf{x}) = y_i - A \sin(\omega t_i + \phi), \quad i = 1, \dots, 21.$$

The Jacobian matrix  $\mathbf{J}(\mathbf{x})$  in this problem is a  $21 \times 3$  matrix with elements given by

$$\begin{aligned} (\mathbf{J}(\mathbf{x}))_{(i,1)} &= -\sin(\omega t_i + \phi), \\ (\mathbf{J}(\mathbf{x}))_{(i,2)} &= -t_i A \cos(\omega t_i + \phi), \\ (\mathbf{J}(\mathbf{x}))_{(i,3)} &= -A \cos(\omega t_i + \phi), \quad i = 1, \dots, 21. \end{aligned}$$

Using the expressions above, we apply the Gauss-Newton algorithm to find the sinusoid of best fit, given the data pairs  $(t_1, y_1), \dots, (t_m, y_m)$ . Figure 9.3 shows a plot of the sinusoid of best fit obtained from the Gauss-Newton algorithm. The parameters of this sinusoid are:  $A = 2.01$ ,  $\omega = 0.992$ , and  $\phi = 0.541$ . ■

A potential problem with the Gauss-Newton method is that the matrix  $\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x})$  may not be positive definite. As described before, this problem can be overcome using a Levenberg-Marquardt modification:

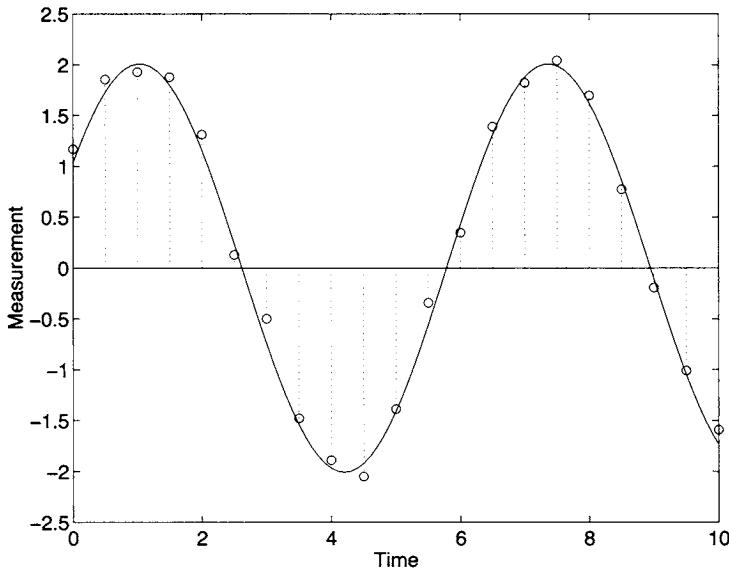
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mu_k \mathbf{I})^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$

This is referred to in the literature as the *Levenberg-Marquardt algorithm*, because the original Levenberg-Marquardt modification was developed specifically for the nonlinear least-squares problem. An alternative interpretation of the Levenberg-Marquardt algorithm is to view the term  $\mu_k \mathbf{I}$  as an approximation to  $\mathbf{S}(\mathbf{x})$  in Newton's algorithm.

## EXERCISES

**9.1** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be given by  $f(x) = (x - x_0)^4$ , where  $x_0 \in \mathbb{R}$  is a constant. Suppose that we apply Newton's method to the problem of minimizing  $f$ .

- a. Write down the update equation for Newton's method applied to the problem.



**Figure 9.3** Sinusoid of best fit in Example 9.3.

- b. Let  $y^{(k)} = |x^{(k)} - x_0|$ , where  $x^{(k)}$  is the  $k$ th iterate in Newton's method. Show that the sequence  $\{y^{(k)}\}$  satisfies  $y^{(k+1)} = \frac{2}{3}y^{(k)}$ .
- c. Show that  $x^{(k)} \rightarrow x_0$  for any initial guess  $x^{(0)}$ .
- d. Show that the order of convergence of the sequence  $\{x^{(k)}\}$  in part b is 1.
- e. Theorem 9.1 states that under certain conditions, the order of convergence of Newton's method is at least 2. Why does that theorem not hold in this particular problem?

**9.2** This question relates to the order of convergence of the secant method, using an argument similar to that of the proof of Theorem 9.1.

- a. Consider a function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^2$ , such that  $x^*$  is a local minimizer and  $f''(x^*) \neq 0$ . Suppose that we apply the algorithm  $x^{(k+1)} = x^{(k)} - \alpha_k f'(x^{(k)})$  such that  $\{\alpha_k\}$  is a positive step-size sequence that converges to  $1/f''(x^*)$ . Show that if  $x^{(k)} \rightarrow x^*$ , then the order of convergence of the algorithm is *superlinear* (i.e., strictly greater than 1).
- b. Given part a, what can you say about the order of convergence of the secant algorithm?

**9.3** Consider the problem of minimizing  $f(x) = x^{\frac{4}{3}} = (\sqrt[3]{x})^4$ ,  $x \in \mathbb{R}$ . Note that 0 is the global minimizer of  $f$ .

- a. Write down the algorithm for Newton's method applied to this problem.
- b. Show that as long as the starting point is not 0, the algorithm in part a does not converge to 0 (no matter how close to 0 we start).

**9.4** Consider *Rosenbrock's Function*:  $f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$ , where  $\mathbf{x} = [x_1, x_2]^\top$  (known to be a “nasty” function—often used as a benchmark for testing algorithms). This function is also known as the *banana function* because of the shape of its level sets.

- a. Prove that  $[1, 1]^\top$  is the unique global minimizer of  $f$  over  $\mathbb{R}^2$ .
- b. With a starting point of  $[0, 0]^\top$ , apply two iterations of Newton's method.  
*Hint:* 
$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$
- c. Repeat part b using a gradient algorithm with a fixed step size of  $\alpha_k = 0.05$  at each iteration.

**9.5** Consider the modified Newton's algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)},$$

where  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)})$ . Suppose that we apply the algorithm to a quadratic function  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Recall that the standard Newton's method reaches point  $\mathbf{x}^*$  such that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  in just one step starting from any initial point  $\mathbf{x}^{(0)}$ . Does the modified Newton's algorithm above possess the same property?



# CHAPTER 10

---

## CONJUGATE DIRECTION METHODS

---

### 10.1 Introduction

The class of *conjugate direction methods* can be viewed as being intermediate between the method of steepest descent and Newton's method. The conjugate direction methods have the following properties:

1. Solve quadratics of  $n$  variables in  $n$  steps.
2. The usual implementation, the *conjugate gradient algorithm*, requires no Hessian matrix evaluations.
3. No matrix inversion and no storage of an  $n \times n$  matrix are required.

The conjugate direction methods typically perform better than the method of steepest descent, but not as well as Newton's method. As we saw from the method of steepest descent and Newton's method, the crucial factor in the efficiency of an iterative search method is the direction of search at each iteration. For a quadratic function of  $n$  variables  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{Q} = \mathbf{Q}^\top > 0$ , the best direction of search, as we shall see, is in the  $\mathbf{Q}$ -conjugate direction. Basically, two directions  $\mathbf{d}^{(1)}$  and  $\mathbf{d}^{(2)}$  in  $\mathbb{R}^n$  are

said to be  $\mathbf{Q}$ -conjugate if  $\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(2)} = 0$ . In general, we have the following definition.

**Definition 10.1** Let  $\mathbf{Q}$  be a real symmetric  $n \times n$  matrix. The directions  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots, \mathbf{d}^{(m)}$  are  $\mathbf{Q}$ -conjugate if for all  $i \neq j$ , we have  $\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$ . ■

**Lemma 10.1** Let  $\mathbf{Q}$  be a symmetric positive definite  $n \times n$  matrix. If the directions  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)} \in \mathbb{R}^n$ ,  $k \leq n - 1$ , are nonzero and  $\mathbf{Q}$ -conjugate, then they are linearly independent. □

*Proof.* Let  $\alpha_0, \dots, \alpha_k$  be scalars such that

$$\alpha_0 \mathbf{d}^{(0)} + \alpha_1 \mathbf{d}^{(1)} + \cdots + \alpha_k \mathbf{d}^{(k)} = \mathbf{0}.$$

Premultiplying this equality by  $\mathbf{d}^{(j)\top} \mathbf{Q}$ ,  $0 \leq j \leq k$ , yields

$$\alpha_j \mathbf{d}^{(j)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0,$$

because all other terms  $\mathbf{d}^{(j)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ ,  $i \neq j$ , by  $\mathbf{Q}$ -conjugacy. But  $\mathbf{Q} = \mathbf{Q}^\top > 0$  and  $\mathbf{d}^{(j)} \neq \mathbf{0}$ ; hence  $\alpha_j = 0$ ,  $j = 0, 1, \dots, k$ . Therefore,  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}$ ,  $k \leq n - 1$ , are linearly independent. ■

**Example 10.1** Let

$$\mathbf{Q} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix}.$$

Note that  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . The matrix  $\mathbf{Q}$  is positive definite because all its leading principal minors are positive:

$$\Delta_1 = 3 > 0, \quad \Delta_2 = \det \begin{bmatrix} 3 & 0 \\ 0 & 4 \end{bmatrix} = 12 > 0, \quad \Delta_3 = \det \mathbf{Q} = 20 > 0.$$

Our goal is to construct a set of  $\mathbf{Q}$ -conjugate vectors  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \mathbf{d}^{(2)}$ .

Let  $\mathbf{d}^{(0)} = [1, 0, 0]^\top$ ,  $\mathbf{d}^{(1)} = [d_1^{(1)}, d_2^{(1)}, d_3^{(1)}]^\top$ ,  $\mathbf{d}^{(2)} = [d_1^{(2)}, d_2^{(2)}, d_3^{(2)}]^\top$ . We require that  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = 0$ . We have

$$\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = [1, 0, 0] \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} d_1^{(1)} \\ d_2^{(1)} \\ d_3^{(1)} \end{bmatrix} = 3d_1^{(1)} + d_3^{(1)}.$$

Let  $d_1^{(1)} = 1$ ,  $d_2^{(1)} = 0$ ,  $d_3^{(1)} = -3$ . Then,  $\mathbf{d}^{(1)} = [1, 0, -3]^\top$ , and thus  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = 0$ .

To find the third vector  $\mathbf{d}^{(2)}$ , which would be  $\mathbf{Q}$ -conjugate with  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$ , we require that  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(2)} = 0$  and  $\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(2)} = 0$ . We have

$$\begin{aligned}\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(2)} &= 3d_1^{(2)} + d_3^{(2)} = 0, \\ \mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(2)} &= -6d_2^{(2)} - 8d_3^{(2)} = 0.\end{aligned}$$

If we take  $\mathbf{d}^{(2)} = [1, 4, -3]^\top$ , then the resulting set of vectors is mutually conjugate. ■

This method of finding  $\mathbf{Q}$ -conjugate vectors is inefficient. A systematic procedure for finding  $\mathbf{Q}$ -conjugate vectors can be devised using the idea underlying the *Gram-Schmidt process* of transforming a given basis of  $\mathbb{R}^n$  into an orthonormal basis of  $\mathbb{R}^n$  (see Exercise 10.1).

## 10.2 The Conjugate Direction Algorithm

We now present the conjugate direction algorithm for minimizing the quadratic function of  $n$  variables

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0, \mathbf{x} \in \mathbb{R}^n$ . Note that because  $\mathbf{Q} > 0$ , the function  $f$  has a global minimizer that can be found by solving  $\mathbf{Q}\mathbf{x} = \mathbf{b}$ .

**Basic Conjugate Direction Algorithm.** Given a starting point  $\mathbf{x}^{(0)}$  and  $\mathbf{Q}$ -conjugate directions  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n-1)}$ ; for  $k \geq 0$ ,

$$\begin{aligned}\mathbf{g}^{(k)} &= \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}, \\ \alpha_k &= -\frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}, \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}.\end{aligned}$$

**Theorem 10.1** *For any starting point  $\mathbf{x}^{(0)}$ , the basic conjugate direction algorithm converges to the unique  $\mathbf{x}^*$  (that solves  $\mathbf{Q}\mathbf{x} = \mathbf{b}$ ) in  $n$  steps; that is,  $\mathbf{x}^{(n)} = \mathbf{x}^*$ .* □

*Proof.* Consider  $\mathbf{x}^* - \mathbf{x}^{(0)} \in \mathbb{R}^n$ . Because the  $\mathbf{d}^{(i)}$  are linearly independent, there exist constants  $\beta_i$ ,  $i = 0, \dots, n-1$ , such that

$$\mathbf{x}^* - \mathbf{x}^{(0)} = \beta_0 \mathbf{d}^{(0)} + \dots + \beta_{n-1} \mathbf{d}^{(n-1)}.$$

Now premultiply both sides of this equation by  $\mathbf{d}^{(k)\top} \mathbf{Q}$ ,  $0 \leq k < n$ , to obtain

$$\mathbf{d}^{(k)\top} \mathbf{Q} (\mathbf{x}^* - \mathbf{x}^{(0)}) = \beta_k \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)},$$

where the terms  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ ,  $k \neq i$ , by the  $\mathbf{Q}$ -conjugate property. Hence,

$$\beta_k = \frac{\mathbf{d}^{(k)\top} \mathbf{Q} (\mathbf{x}^* - \mathbf{x}^{(0)})}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

Now, we can write

$$\mathbf{x}^{(k)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} + \cdots + \alpha_{k-1} \mathbf{d}^{(k-1)}.$$

Therefore,

$$\mathbf{x}^{(k)} - \mathbf{x}^{(0)} = \alpha_0 \mathbf{d}^{(0)} + \cdots + \alpha_{k-1} \mathbf{d}^{(k-1)}.$$

So writing

$$\mathbf{x}^* - \mathbf{x}^{(0)} = (\mathbf{x}^* - \mathbf{x}^{(k)}) + (\mathbf{x}^{(k)} - \mathbf{x}^{(0)})$$

and premultiplying the above by  $\mathbf{d}^{(k)\top} \mathbf{Q}$ , we obtain

$$\mathbf{d}^{(k)\top} \mathbf{Q} (\mathbf{x}^* - \mathbf{x}^{(0)}) = \mathbf{d}^{(k)\top} \mathbf{Q} (\mathbf{x}^* - \mathbf{x}^{(k)}) = -\mathbf{d}^{(k)\top} \mathbf{g}^{(k)},$$

because  $\mathbf{g}^{(k)} = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}$  and  $\mathbf{Q} \mathbf{x}^* = \mathbf{b}$ . Thus,

$$\beta_k = -\frac{\mathbf{d}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}} = \alpha_k$$

and  $\mathbf{x}^* = \mathbf{x}^{(n)}$ , which completes the proof. ■

**Example 10.2** Find the minimizer of

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \mathbf{x} \in \mathbb{R}^2,$$

using the conjugate direction method with the initial point  $\mathbf{x}^{(0)} = [0, 0]^\top$ , and  $\mathbf{Q}$ -conjugate directions  $\mathbf{d}^{(0)} = [1, 0]^\top$  and  $\mathbf{d}^{(1)} = [-\frac{3}{8}, \frac{3}{4}]^\top$ .

We have

$$\mathbf{g}^{(0)} = -\mathbf{b} = [1, -1]^\top,$$

and hence

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = -\frac{[1, -1]^\top \begin{bmatrix} 1 \\ 0 \end{bmatrix}}{[1, 0]^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}} = -\frac{1}{4}.$$

Thus,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \frac{1}{4} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -\frac{1}{4} \\ 0 \end{bmatrix}.$$

To find  $\mathbf{x}^{(2)}$ , we compute

$$\mathbf{g}^{(1)} = \mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -\frac{1}{4} \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ -\frac{3}{2} \end{bmatrix}$$

and

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = -\frac{\begin{bmatrix} 0, -\frac{3}{2} \end{bmatrix} \begin{bmatrix} -\frac{3}{8} \\ \frac{3}{4} \end{bmatrix}}{\begin{bmatrix} -\frac{3}{8}, \frac{3}{4} \end{bmatrix} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -\frac{3}{8} \\ \frac{3}{4} \end{bmatrix}} = 2.$$

Therefore,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = \begin{bmatrix} -\frac{1}{4} \\ 0 \end{bmatrix} + 2 \begin{bmatrix} -\frac{3}{8} \\ \frac{3}{4} \end{bmatrix} = \begin{bmatrix} -1 \\ \frac{3}{2} \end{bmatrix}.$$

Because  $f$  is a quadratic function in two variables,  $\mathbf{x}^{(2)} = \mathbf{x}^*$ . ■

For a quadratic function of  $n$  variables, the conjugate direction method reaches the solution after  $n$  steps. As we shall see below, the method also possesses a certain desirable property in the intermediate steps. To see this, suppose that we start at  $\mathbf{x}^{(0)}$  and search in the direction  $\mathbf{d}^{(0)}$  to obtain

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \left( \frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} \right) \mathbf{d}^{(0)}.$$

We claim that

$$\mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0.$$

To see this,

$$\begin{aligned} \mathbf{g}^{(1)\top} \mathbf{d}^{(0)} &= (\mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b})^\top \mathbf{d}^{(0)} \\ &= \mathbf{x}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)} - \left( \frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} \right) \mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)} - \mathbf{b}^\top \mathbf{d}^{(0)} \\ &= \mathbf{g}^{(0)\top} \mathbf{d}^{(0)} - \mathbf{g}^{(0)\top} \mathbf{d}^{(0)} = 0. \end{aligned}$$

The equation  $\mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0$  implies that  $\alpha_0$  has the property that  $\alpha_0 = \arg \min \phi_0(\alpha)$ , where  $\phi_0(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$ . To see this, apply the chain rule to get

$$\frac{d\phi_0}{d\alpha}(\alpha) = \nabla f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})^\top \mathbf{d}^{(0)}.$$

Evaluating the above at  $\alpha = \alpha_0$ , we get

$$\frac{d\phi_0}{d\alpha}(\alpha_0) = \mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0.$$

Because  $\phi_0$  is a quadratic function of  $\alpha$ , and the coefficient of the  $\alpha^2$  term in  $\phi_0$  is  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)} > 0$ , the above implies that  $\alpha_0 = \arg \min_{\alpha \in \mathbb{R}} \phi_0(\alpha)$ .

Using a similar argument, we can show that for all  $k$ ,

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} = 0$$

and hence

$$\alpha_k = \arg \min f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

In fact, an even stronger condition holds, as given by the following lemma.

**Lemma 10.2** *In the conjugate direction algorithm,*

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0$$

for all  $k$ ,  $0 \leq k \leq n - 1$ , and  $0 \leq i \leq k$ .  $\square$

*Proof.* Note that

$$\mathbf{Q}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \mathbf{Q}\mathbf{x}^{(k+1)} - \mathbf{b} - (\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}) = \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)},$$

because  $\mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}$ . Thus,

$$\mathbf{g}^{(k+1)} = \mathbf{g}^{(k)} + \alpha_k \mathbf{Q} \mathbf{d}^{(k)}.$$

We prove the lemma by induction. The result is true for  $k = 0$  because  $\mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0$ , as shown before. We now show that if the result is true for  $k - 1$  (i.e.,  $\mathbf{g}^{(k)\top} \mathbf{d}^{(i)} = 0$ ,  $i \leq k - 1$ ), then it is true for  $k$  (i.e.,  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0$ ,  $i \leq k$ ). Fix  $k > 0$  and  $0 \leq i < k$ . By the induction hypothesis,  $\mathbf{g}^{(k)\top} \mathbf{d}^{(i)} = 0$ . Because

$$\mathbf{g}^{(k+1)} = \mathbf{g}^{(k)} + \alpha_k \mathbf{Q} \mathbf{d}^{(k)},$$

and  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$  by  $\mathbf{Q}$ -conjugacy, we have

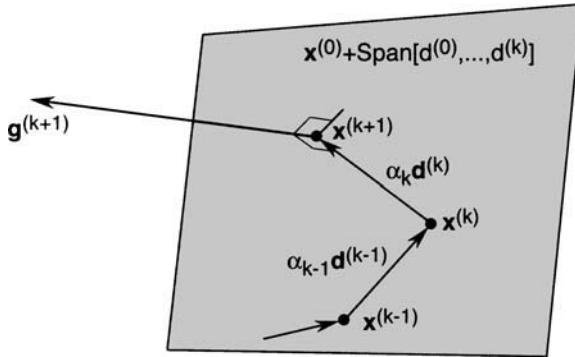
$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = \mathbf{g}^{(k)\top} \mathbf{d}^{(i)} + \alpha_k \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0.$$

It remains to be shown that

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} = 0.$$

Indeed,

$$\begin{aligned} \mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} &= (\mathbf{Q}\mathbf{x}^{(k+1)} - \mathbf{b})^\top \mathbf{d}^{(k)} \\ &= \left( \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}} \mathbf{d}^{(k)} \right)^\top \mathbf{Q} \mathbf{d}^{(k)} - \mathbf{b}^\top \mathbf{d}^{(k)} \\ &= (\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b})^\top \mathbf{d}^{(k)} - \mathbf{g}^{(k)\top} \mathbf{d}^{(k)} \\ &= 0, \end{aligned}$$



**Figure 10.1** Illustration of Lemma 10.2.

because  $\mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b} = \mathbf{g}^{(k)}$ .

Therefore, by induction, for all  $0 \leq k \leq n - 1$  and  $0 \leq i \leq k$ ,

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0.$$

■

By Lemma 10.2 we see that  $\mathbf{g}^{(k+1)}$  is orthogonal to any vector from the subspace spanned by  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}$ . Figure 10.1 illustrates this statement.

The lemma can be used to show an interesting optimal property of the conjugate direction algorithm. Specifically, we now show that not only does  $f(\mathbf{x}^{(k+1)})$  satisfy  $f(\mathbf{x}^{(k+1)}) = \min_{\alpha} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ , as indicated before, but also

$$f(\mathbf{x}^{(k+1)}) = \min_{a_0, \dots, a_k} f\left(\mathbf{x}^{(0)} + \sum_{i=0}^k a_i \mathbf{d}^{(i)}\right).$$

In other words, if we write

$$\mathcal{V}_k = \mathbf{x}^{(0)} + \text{span}[\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}],$$

then we can express  $f(\mathbf{x}^{(k+1)}) = \min_{\mathbf{x} \in \mathcal{V}_k} f(\mathbf{x})$ . As  $k$  increases, the subspace  $\text{span}[\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}]$  “expands,” and will eventually fill the whole of  $\mathbb{R}^n$  (provided that the vectors  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots$  are linearly independent). Therefore, for some sufficiently large  $k$ ,  $\mathbf{x}^*$  will lie in  $\mathcal{V}_k$ . For this reason, the above result is sometimes called the *expanding subspace theorem* (see, e.g., [88, p. 266]).

To prove the expanding subspace theorem, define the matrix  $\mathbf{D}^{(k)}$  by

$$\mathbf{D}^{(k)} = [\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}];$$

that is,  $\mathbf{d}^{(i)}$  is the  $i$ th column of  $\mathbf{D}^{(k)}$ . Note that  $\mathbf{x}^{(0)} + \mathcal{R}(\mathbf{D}^{(k)}) = \mathcal{V}_k$ . Also,

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \mathbf{x}^{(0)} + \sum_{i=0}^k \alpha_i \mathbf{d}^{(i)} \\ &= \mathbf{x}^{(0)} + \mathbf{D}^{(k)} \boldsymbol{\alpha},\end{aligned}$$

where  $\boldsymbol{\alpha} = [\alpha_0, \dots, \alpha_k]^\top$ . Hence,

$$\mathbf{x}^{(k+1)} \in \mathbf{x}^{(0)} + \mathcal{R}(\mathbf{D}^{(k)}) = \mathcal{V}_k.$$

Now, consider any vector  $\mathbf{x} \in \mathcal{V}_k$ . There exists a vector  $\mathbf{a}$  such that  $\mathbf{x} = \mathbf{x}^{(0)} + \mathbf{D}^{(k)} \mathbf{a}$ . Let  $\phi_k(\mathbf{a}) = f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)} \mathbf{a})$ . Note that  $\phi_k$  is a quadratic function and has a unique minimizer that satisfies the FONC (see Exercises 6.33 and 10.7). By the chain rule,

$$D\phi_k(\mathbf{a}) = \nabla f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)} \mathbf{a})^\top \mathbf{D}^{(k)}.$$

Therefore,

$$\begin{aligned}D\phi_k(\boldsymbol{\alpha}) &= \nabla f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)} \boldsymbol{\alpha})^\top \mathbf{D}^{(k)} \\ &= \nabla f(\mathbf{x}^{(k+1)})^\top \mathbf{D}^{(k)} \\ &= \mathbf{g}^{(k+1)^\top} \mathbf{D}^{(k)}.\end{aligned}$$

By Lemma 10.2,  $\mathbf{g}^{(k+1)^\top} \mathbf{D}^{(k)} = \mathbf{0}^\top$ . Therefore,  $\boldsymbol{\alpha}$  satisfies the FONC for the quadratic function  $\phi_k$ , and hence  $\boldsymbol{\alpha}$  is the minimizer of  $\phi_k$ ; that is,

$$f(\mathbf{x}^{(k+1)}) = \min_{\mathbf{a}} f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)} \mathbf{a}) = \min_{\mathbf{x} \in \mathcal{V}_k} f(\mathbf{x}),$$

which completes the proof of our result.

The conjugate direction algorithm is very effective. However, to use the algorithm, we need to specify the  $\mathbf{Q}$ -conjugate directions. Fortunately, there is a way to generate  $\mathbf{Q}$ -conjugate directions as we perform iterations. In the next section we discuss an algorithm that incorporates the generation of  $\mathbf{Q}$ -conjugate directions.

### 10.3 The Conjugate Gradient Algorithm

The conjugate gradient algorithm does not use prespecified conjugate directions, but instead computes the directions as the algorithm progresses. At each stage of the algorithm, the direction is calculated as a linear combination of the previous direction and the current gradient, in such a way that all the directions are mutually  $\mathbf{Q}$ -conjugate—hence the name *conjugate gradient algorithm*. This calculation exploits the fact that for a quadratic function of

$n$  variables, we can locate the function minimizer by performing  $n$  searches along mutually conjugate directions.

As before, we consider the quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b}, \quad \mathbf{x} \in \mathbb{R}^n,$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Our first search direction from an initial point  $\mathbf{x}^{(0)}$  is in the direction of steepest descent; that is,

$$\mathbf{d}^{(0)} = -\mathbf{g}^{(0)}.$$

Thus,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)},$$

where

$$\alpha_0 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}) = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}}.$$

In the next stage, we search in a direction  $\mathbf{d}^{(1)}$  that is  $\mathbf{Q}$ -conjugate to  $\mathbf{d}^{(0)}$ . We choose  $\mathbf{d}^{(1)}$  as a linear combination of  $\mathbf{g}^{(1)}$  and  $\mathbf{d}^{(0)}$ . In general, at the  $(k+1)$ th step, we choose  $\mathbf{d}^{(k+1)}$  to be a linear combination of  $\mathbf{g}^{(k+1)}$  and  $\mathbf{d}^{(k)}$ . Specifically, we choose

$$\mathbf{d}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)}, \quad k = 0, 1, 2, \dots$$

The coefficients  $\beta_k$ ,  $k = 1, 2, \dots$ , are chosen in such a way that  $\mathbf{d}^{(k+1)}$  is  $\mathbf{Q}$ -conjugate to  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}$ . This is accomplished by choosing  $\beta_k$  to be

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

The conjugate gradient algorithm is summarized below.

1. Set  $k := 0$ ; select the initial point  $\mathbf{x}^{(0)}$ .
2.  $\mathbf{g}^{(0)} = \nabla f(\mathbf{x}^{(0)})$ . If  $\mathbf{g}^{(0)} = \mathbf{0}$ , stop; else, set  $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)}$ .
3.  $\alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}$ .
4.  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ .
5.  $\mathbf{g}^{(k+1)} = \nabla f(\mathbf{x}^{(k+1)})$ . If  $\mathbf{g}^{(k+1)} = \mathbf{0}$ , stop.
6.  $\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}$ .
7.  $\mathbf{d}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)}$ .
8. Set  $k := k + 1$ ; go to step 3.

**Proposition 10.1** *In the conjugate gradient algorithm, the directions  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n-1)}$  are  $\mathbf{Q}$ -conjugate.*  $\square$

*Proof.* We use induction. We first show that  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = 0$ . To this end we write

$$\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = \mathbf{d}^{(0)\top} \mathbf{Q}(-\mathbf{g}^{(1)} + \beta_0 \mathbf{d}^{(0)}).$$

Substituting for

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}}$$

in the equation above, we see that  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = 0$ .

We now assume that  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}$ ,  $k < n - 1$ , are  $\mathbf{Q}$ -conjugate directions. From Lemma 10.2 we have  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(j)} = 0$ ,  $j = 0, 1, \dots, k$ . Thus,  $\mathbf{g}^{(k+1)}$  is orthogonal to each of the directions  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots, \mathbf{d}^{(k)}$ . We now show that

$$\mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} = 0, \quad j = 0, 1, \dots, k.$$

Fix  $j \in \{0, \dots, k\}$ . We have

$$\mathbf{d}^{(j)} = -\mathbf{g}^{(j)} + \beta_{j-1} \mathbf{d}^{(j-1)}.$$

Substituting this equation into the previous one yields

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(j)} = 0 = -\mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} + \beta_{j-1} \mathbf{g}^{(k+1)\top} \mathbf{d}^{(j-1)}.$$

Because  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(j-1)} = 0$ , it follows that  $\mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} = 0$ .

We are now ready to show that  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$ ,  $j = 0, \dots, k$ . We have

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = (-\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)})^\top \mathbf{Q} \mathbf{d}^{(j)}.$$

If  $j < k$ , then  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0$ , by virtue of the induction hypothesis. Hence, we have

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = -\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)}.$$

But  $\mathbf{g}^{(j+1)} = \mathbf{g}^{(j)} + \alpha_j \mathbf{Q} \mathbf{d}^{(j)}$ . Because  $\mathbf{g}^{(k+1)\top} \mathbf{g}^{(i)} = 0$ ,  $i = 0, \dots, k$ ,

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = -\mathbf{g}^{(k+1)\top} \frac{(\mathbf{g}^{(j+1)} - \mathbf{g}^{(j)})}{\alpha_j} = 0.$$

Thus,

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(j)} = 0, \quad j = 0, \dots, k-1.$$

It remains to be shown that  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)} = 0$ . We have

$$\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)} = (-\mathbf{g}^{(k+1)} + \beta_k \mathbf{d}^{(k)})^\top \mathbf{Q} \mathbf{d}^{(k)}.$$

Using the expression for  $\beta_k$ , we get  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)} = 0$ , which completes the proof.  $\blacksquare$

**Example 10.3** Consider the quadratic function

$$f(x_1, x_2, x_3) = \frac{3}{2}x_1^2 + 2x_2^2 + \frac{3}{2}x_3^2 + x_1x_3 + 2x_2x_3 - 3x_1 - x_3.$$

We find the minimizer using the conjugate gradient algorithm, using the starting point  $\mathbf{x}^{(0)} = [0, 0, 0]^\top$ .

We can represent  $f$  as

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where

$$\mathbf{Q} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}.$$

We have

$$\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b} = [3x_1 + x_3 - 3, 4x_2 + 2x_3, x_1 + 2x_2 + 3x_3 - 1]^\top.$$

Hence,

$$\begin{aligned} \mathbf{g}^{(0)} &= [-3, 0, -1]^\top, \\ \mathbf{d}^{(0)} &= -\mathbf{g}^{(0)}, \\ \alpha_0 &= -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{10}{36} = 0.2778 \end{aligned}$$

and

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [0.8333, 0, 0.2778]^\top.$$

The next stage yields

$$\begin{aligned} \mathbf{g}^{(1)} &= \nabla f(\mathbf{x}^{(1)}) = [-0.2222, 0.5556, 0.6667]^\top, \\ \beta_0 &= \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = 0.08025. \end{aligned}$$

We can now compute

$$\mathbf{d}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{d}^{(0)} = [0.4630, -0.5556, -0.5864]^\top.$$

Hence,

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 0.2187$$

and

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [0.9346, -0.1215, 0.1495]^\top.$$

To perform the third iteration, we compute

$$\begin{aligned}\mathbf{g}^{(2)} &= \nabla f(\mathbf{x}^{(2)}) = [-0.04673, -0.1869, 0.1402]^\top, \\ \beta_1 &= \frac{\mathbf{g}^{(2)\top} \mathbf{Q} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 0.07075, \\ \mathbf{d}^{(2)} &= -\mathbf{g}^{(2)} + \beta_1 \mathbf{d}^{(1)} = [0.07948, 0.1476, -0.1817]^\top.\end{aligned}$$

Hence,

$$\alpha_2 = -\frac{\mathbf{g}^{(2)\top} \mathbf{d}^{(2)}}{\mathbf{d}^{(2)\top} \mathbf{Q} \mathbf{d}^{(2)}} = 0.8231$$

and

$$\mathbf{x}^{(3)} = \mathbf{x}^{(2)} + \alpha_2 \mathbf{d}^{(2)} = [1.000, 0.000, 0.000]^\top.$$

Note that

$$\mathbf{g}^{(3)} = \nabla f(\mathbf{x}^{(3)}) = \mathbf{0},$$

as expected, because  $f$  is a quadratic function of three variables. Hence,  $\mathbf{x}^* = \mathbf{x}^{(3)}$ . ■

## 10.4 The Conjugate Gradient Algorithm for Nonquadratic Problems

In Section 10.3, we showed that the conjugate gradient algorithm is a conjugate direction method, and therefore minimizes a positive definite quadratic function of  $n$  variables in  $n$  steps. The algorithm can be extended to general nonlinear functions by interpreting  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$  as a second-order Taylor series approximation of the objective function. Near the solution such functions behave approximately as quadratics, as suggested by the Taylor series expansion. For a quadratic, the matrix  $\mathbf{Q}$ , the Hessian of the quadratic, is constant. However, for a general nonlinear function the Hessian is a matrix that has to be reevaluated at each iteration of the algorithm. This can be computationally very expensive. Thus, an efficient implementation of the conjugate gradient algorithm that eliminates the Hessian evaluation at each step is desirable.

Observe that  $\mathbf{Q}$  appears only in the computation of the scalars  $\alpha_k$  and  $\beta_k$ . Because

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}),$$

the closed-form formula for  $\alpha_k$  in the algorithm can be replaced by a numerical line search procedure. Therefore, we need only concern ourselves with the formula for  $\beta_k$ . Fortunately, elimination of  $\mathbf{Q}$  from the formula is possible and results in algorithms that depend only on the function and gradient values at

each iteration. We now discuss modifications of the conjugate gradient algorithm for a quadratic function for the case in which the Hessian is unknown but in which objective function values and gradients are available. The modifications are all based on algebraically manipulating the formula  $\beta_k$  in such a way that  $\mathbf{Q}$  is eliminated. We discuss three well-known modifications.

**Hestenes-Stiefel Formula.** Recall that

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)}}.$$

The Hestenes-Stiefel formula is based on replacing the term  $\mathbf{Q} \mathbf{d}^{(k)}$  by the term  $(\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)})/\alpha_k$ . The two terms are equal in the quadratic case, as we now show. Now,  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ . Premultiplying both sides by  $\mathbf{Q}$ , subtracting  $\mathbf{b}$  from both sides, and recognizing that  $\mathbf{g}^{(k)} = \mathbf{Q} \mathbf{x}^{(k)} - \mathbf{b}$ , we get  $\mathbf{g}^{(k+1)} = \mathbf{g}^{(k)} + \alpha_k \mathbf{Q} \mathbf{d}^{(k)}$ , which we can rewrite as  $\mathbf{Q} \mathbf{d}^{(k)} = (\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)})/\alpha_k$ . Substituting this into the original equation for  $\beta_k$  gives the *Hestenes-Stiefel formula*

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{d}^{(k)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}.$$

**Polak-Ribière Formula.** Starting from the Hestenes-Stiefel formula, we multiply out the denominator to get

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} - \mathbf{d}^{(k)\top} \mathbf{g}^{(k)}}.$$

By Lemma 10.2,  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} = 0$ . Also, since  $\mathbf{d}^{(k)} = -\mathbf{g}^{(k)} + \beta_{k-1} \mathbf{d}^{(k-1)}$ , and premultiplying this by  $\mathbf{g}^{(k)\top}$ , we get

$$\mathbf{g}^{(k)\top} \mathbf{d}^{(k)} = -\mathbf{g}^{(k)\top} \mathbf{g}^{(k)} + \beta_{k-1} \mathbf{g}^{(k)\top} \mathbf{d}^{(k-1)} = -\mathbf{g}^{(k)\top} \mathbf{g}^{(k)},$$

where once again we used Lemma 10.2. Hence, we get the Polak-Ribière formula

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}.$$

**Fletcher-Reeves Formula.** Starting with the Polak-Ribière formula, we multiply out the numerator to get

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k+1)} - \mathbf{g}^{(k+1)\top} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}.$$

We now use the fact that  $\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k)} = 0$ , which we get by using the equation

$$\mathbf{g}^{(k+1)\top} \mathbf{d}^{(k)} = -\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k)} + \beta_{k-1} \mathbf{g}^{(k+1)\top} \mathbf{d}^{(k-1)}$$

and applying Lemma 10.2. This leads to the Fletcher-Reeves formula

$$\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k+1)}}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}}.$$

The formulas above give us conjugate gradient algorithms that do not require explicit knowledge of the Hessian matrix  $\mathbf{Q}$ . All we need are the objective function and gradient values at each iteration. For the quadratic case the three expressions for  $\beta_k$  are exactly equal. However, this is not the case for a general nonlinear objective function.

We need a few more slight modifications to apply the algorithm to general nonlinear functions in practice. First, as mentioned in our discussion of the steepest descent algorithm (Section 8.2), the stopping criterion  $\nabla f(\mathbf{x}^{(k+1)}) = \mathbf{0}$  is not practical. A suitable practical stopping criterion, such as those discussed in Section 8.2, needs to be used.

For nonquadratic problems, the algorithm will not usually converge in  $n$  steps, and as the algorithm progresses, the “ $\mathbf{Q}$ -conjugacy” of the direction vectors will tend to deteriorate. Thus, a common practice is to reinitialize the direction vector to the negative gradient after every few iterations (e.g.,  $n$  or  $n + 1$ ) and continue until the algorithm satisfies the stopping criterion.

A very important issue in minimization problems of nonquadratic functions is the line search. The purpose of the line search is to minimize  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$  with respect to  $\alpha \geq 0$ . A typical approach is to bracket or box in the minimizer and then estimate it. The accuracy of the line search is a critical factor in the performance of the conjugate gradient algorithm. If the line search is known to be inaccurate, the Hestenes-Stiefel formula for  $\beta_k$  is recommended [69].

In general, the choice of which formula for  $\beta_k$  to use depends on the objective function. For example, the Polak-Ribière formula is known to perform far better than the Fletcher-Reeves formula in some cases but not in others. In fact, there are cases in which the  $\mathbf{g}^{(k)}$ ,  $k = 1, 2, \dots$ , are bounded away from zero when the Polak-Ribière formula is used (see [107]). In the study by Powell in [107], a global convergence analysis suggests that the Fletcher-Reeves formula for  $\beta_k$  is superior. Powell further suggests another formula for  $\beta_k$ :

$$\beta_k = \max \left\{ 0, \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}} \right\}.$$

For general results on the convergence of conjugate gradient methods, we refer the reader to [135]. For an application of conjugate gradient algorithms to Wiener filtering, see [116], [117], and [118].

Conjugate gradient algorithms are related to *Krylov subspace methods* (see Exercise 10.6). Krylov-subspace-iteration methods, initiated by Magnus Hestenes, Eduard Stiefel, and Cornelius Lanczos, have been declared one of the 10 algorithms with the greatest influence on the development and practice of science and engineering in the twentieth century [40].

For control perspective on the conjugate gradient algorithm, derived from a proportional-plus-derivative (PD) controller architecture, see [4]. In addition, these authors offer a control perspective on Krylov-subspace-iteration methods as discrete feedback control systems.

## EXERCISES

**10.1** (Adopted from [88, Exercise 9.8(1)]) Let  $\mathbf{Q}$  be a real symmetric positive definite  $n \times n$  matrix. Given an arbitrary set of linearly independent vectors  $\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(n-1)}\}$  in  $\mathbb{R}^n$ , the *Gram-Schmidt procedure* generates a set of vectors  $\{\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(n-1)}\}$  as follows:

$$\begin{aligned}\mathbf{d}^{(0)} &= \mathbf{p}^{(0)}, \\ \mathbf{d}^{(k+1)} &= \mathbf{p}^{(k+1)} - \sum_{i=0}^k \frac{\mathbf{p}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)}}{\mathbf{d}^{(i)\top} \mathbf{Q} \mathbf{d}^{(i)}} \mathbf{d}^{(i)}.\end{aligned}$$

Show that the vectors  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(n-1)}$  are  $\mathbf{Q}$ -conjugate.

**10.2** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be the quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Given a set of directions  $\{\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots\} \subset \mathbb{R}^n$ , consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

where  $\alpha_k$  is the step size. Suppose that  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0$  for all  $k = 0, \dots, n-1$  and  $i = 0, \dots, k$ , where  $\mathbf{g}^{(k+1)} = \nabla f(\mathbf{x}^{(k+1)})$ . Show that if  $\mathbf{g}^{(k)\top} \mathbf{d}^{(k)} \neq 0$  for all  $k = 0, \dots, n-1$ , then  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(n-1)}$  are  $\mathbf{Q}$ -conjugate.

**10.3** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , where  $\mathbf{b} \in \mathbb{R}^n$  and  $\mathbf{Q}$  is a real symmetric positive definite  $n \times n$  matrix. Show that in the conjugate gradient method for this  $f$ ,  $\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k)} = -\mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{g}^{(k)}$ .

**10.4** Let  $\mathbf{Q}$  be a real  $n \times n$  symmetric matrix.

- a. Show that there exists a  $\mathbf{Q}$ -conjugate set  $\{\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n)}\}$  such that each  $\mathbf{d}^{(i)}$  ( $i = 1, \dots, n$ ) is an eigenvector of  $\mathbf{Q}$ .  
*Hint:* Use the fact that for any real symmetric  $n \times n$  matrix, there exists a set  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  of its eigenvectors such that  $\mathbf{v}_i^\top \mathbf{v}_j = 0$  for all  $i, j = 1, \dots, n$ ,  $i \neq j$ .
- b. Suppose that  $\mathbf{Q}$  is positive definite. Show that if  $\{\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(n)}\}$  is a  $\mathbf{Q}$ -conjugate set that is also orthogonal (i.e.,  $\mathbf{d}^{(i)\top} \mathbf{d}^{(j)} = 0$  for all  $i, j = 1, \dots, n$ ,  $i \neq j$ ), and  $\mathbf{d}^{(i)} \neq \mathbf{0}$ ,  $i = 1, \dots, n$ , then each  $\mathbf{d}^{(i)}$ ,  $i = 1, \dots, n$ , is an eigenvector of  $\mathbf{Q}$ .

**10.5** Consider the following algorithm for minimizing a function  $f$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

where  $\alpha_k = \arg \min_{\alpha} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . Let  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$  (as usual).

Suppose that  $f$  is quadratic with Hessian  $\mathbf{Q}$ . We choose  $\mathbf{d}^{(k+1)} = \gamma_k \mathbf{g}^{(k+1)} + \mathbf{d}^{(k)}$ , and we wish the directions  $\mathbf{d}^{(k)}$  and  $\mathbf{d}^{(k+1)}$  to be  $\mathbf{Q}$ -conjugate. Find a formula for  $\gamma_k$  in terms of  $\mathbf{d}^{(k)}$ ,  $\mathbf{g}^{(k+1)}$ , and  $\mathbf{Q}$ .

**10.6** Consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

with  $\alpha_k \in \mathbb{R}$  scalar and  $\mathbf{x}^{(0)} = \mathbf{0}$ , applied to the quadratic function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

where  $\mathbf{Q} > 0$ . As usual, write  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ . Suppose that the search directions are generated according to

$$\mathbf{d}^{(k+1)} = a_k \mathbf{g}^{(k+1)} + b_k \mathbf{d}^{(k)},$$

where  $a_k$  and  $b_k$  are real constants, and by convention we take  $\mathbf{d}^{(-1)} = \mathbf{0}$ .

- a. Define the subspace  $\mathcal{V}_k = \text{span}[\mathbf{b}, \mathbf{Q}\mathbf{b}, \dots, \mathbf{Q}^{k-1}\mathbf{b}]$  (called the *Krylov subspace of order k*). Show that  $\mathbf{d}^{(k)} \in \mathcal{V}_{k+1}$  and  $\mathbf{x}^{(k)} \in \mathcal{V}_k$ .

*Hint:* Use induction. Note that  $\mathcal{V}_0 = \{\mathbf{0}\}$  and  $\mathcal{V}_1 = \text{span}[\mathbf{b}]$ .

- b. In light of part a, what can you say about the “optimality” of the conjugate gradient algorithm with respect to the Krylov subspace?

**10.7** Consider the quadratic function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Let  $\mathbf{D} \in \mathbb{R}^{n \times r}$  be of rank  $r$  and  $\mathbf{x}_0 \in \mathbb{R}^n$ . Define the function  $\phi : \mathbb{R}^r \rightarrow \mathbb{R}$  by

$$\phi(\mathbf{a}) = f(\mathbf{x}_0 + \mathbf{D}\mathbf{a}).$$

Show that  $\phi$  is a quadratic function with a positive definite quadratic term.

**10.8** Consider a conjugate gradient algorithm applied to a quadratic function.

- a. Show that the gradients associated with the algorithm are mutually orthogonal. Specifically, show that  $\mathbf{g}^{(k+1)\top} \mathbf{g}^{(i)} = 0$  for all  $0 \leq k \leq n-1$  and  $0 \leq i \leq k$ .

*Hint:* Write  $\mathbf{g}^{(i)}$  in terms of  $\mathbf{d}^{(i)}$  and  $\mathbf{d}^{(i-1)}$ .

- b. Show that the gradients associated with the algorithm are  $\mathbf{Q}$ -conjugate if separated by at least two iterations. Specifically, show that  $\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{g}^{(i)} = 0$  for all  $0 \leq k \leq n-1$  and  $0 \leq i \leq k-1$ .

**10.9** Represent the function

$$f(x_1, x_2) = \frac{5}{2}x_1^2 + x_2^2 - 3x_1x_2 - x_2 - 7$$

in the form  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b} + c$ . Then use the *conjugate gradient algorithm* to construct a vector  $\mathbf{d}^{(1)}$  that is  $\mathbf{Q}$ -conjugate with  $\mathbf{d}^{(0)} = \nabla f(\mathbf{x}^{(0)})$ , where  $\mathbf{x}^{(0)} = \mathbf{0}$ .

**10.10** Let  $f(\mathbf{x})$ ,  $\mathbf{x} = [x_1, x_2]^\top \in \mathbb{R}^2$ , be given by

$$f(\mathbf{x}) = \frac{5}{2}x_1^2 + \frac{1}{2}x_2^2 + 2x_1x_2 - 3x_1 - x_2.$$

- a. Express  $f(\mathbf{x})$  in the form of  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ .
- b. Find the minimizer of  $f$  using the conjugate gradient algorithm. Use a starting point of  $\mathbf{x}^{(0)} = [0, 0]^\top$ .
- c. Calculate the minimizer of  $f$  analytically from  $\mathbf{Q}$  and  $\mathbf{b}$ , and check it with your answer in part b.

**10.11** Write a MATLAB program to implement the conjugate gradient algorithm for general functions. Use the secant method for the line search (e.g., the MATLAB function of Exercise 7.11). Test the different formulas for  $\beta_k$  on Rosenbrock's function (see Exercise 9.4) with an initial condition  $\mathbf{x}^{(0)} = [-2, 2]^\top$ . For this exercise, reinitialize the update direction to the negative gradient every six iterations.



# CHAPTER 11

---

## QUASI-NEWTON METHODS

---

### 11.1 Introduction

Newton's method is one of the more successful algorithms for optimization. If it converges, it has a quadratic order of convergence. However, as pointed out before, for a general nonlinear objective function, convergence to a solution cannot be guaranteed from an arbitrary initial point  $\mathbf{x}^{(0)}$ . In general, if the initial point is not sufficiently close to the solution, then the algorithm may not possess the descent property [i.e.,  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$  for some  $k$ ].

Recall that the idea behind Newton's method is to locally approximate the function  $f$  being minimized, at every iteration, by a quadratic function. The minimizer for the quadratic approximation is used as the starting point for the next iteration. This leads to Newton's recursive algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}.$$

We may try to guarantee that the algorithm has the descent property by modifying the original algorithm as follows:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)},$$

where  $\alpha_k$  is chosen to ensure that

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}).$$

For example, we may choose  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)})$  (see Theorem 9.2). We can then determine an appropriate value of  $\alpha_k$  by performing a line search in the direction  $-\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ . Note that although the line search is simply the minimization of the real variable function  $\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)})$ , it is not a trivial problem to solve.

A computational drawback of Newton's method is the need to evaluate  $\mathbf{F}(\mathbf{x}^{(k)})$  and solve the equation  $\mathbf{F}(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$  [i.e., compute  $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ ]. To avoid the computation of  $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$ , the quasi-Newton methods use an approximation to  $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$  in place of the true inverse. This approximation is updated at every stage so that it exhibits at least some properties of  $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$ . To get some idea about the properties that an approximation to  $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$  should satisfy, consider the formula

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)},$$

where  $\mathbf{H}_k$  is an  $n \times n$  real matrix and  $\alpha$  is a positive search parameter. Expanding  $f$  about  $\mathbf{x}^{(k)}$  yields

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k)}) + \mathbf{g}^{(k)\top} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + o(\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|) \\ &= f(\mathbf{x}^{(k)}) - \alpha \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)} + o(\|\mathbf{H}_k \mathbf{g}^{(k)}\| \alpha). \end{aligned}$$

As  $\alpha$  tends to zero, the second term on the right-hand side of this equation dominates the third. Thus, to guarantee a decrease in  $f$  for small  $\alpha$ , we have to have

$$\mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)} > 0.$$

A simple way to ensure this is to require that  $\mathbf{H}_k$  be positive definite. We have proved the following result.

**Proposition 11.1** *Let  $f \in \mathcal{C}^1$ ,  $\mathbf{x}^{(k)} \in \mathbb{R}^n$ ,  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , and  $\mathbf{H}_k$  an  $n \times n$  real symmetric positive definite matrix. If we set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k \mathbf{g}^{(k)}$ , where  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)})$ , then  $\alpha_k > 0$  and  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ .*  $\square$

In constructing an approximation to the inverse of the Hessian matrix, we should use only the objective function and gradient values. Thus, if we can find a suitable method of choosing  $\mathbf{H}_k$ , the iteration may be carried out without any evaluation of the Hessian and without the solution of any set of linear equations.

## 11.2 Approximating the Inverse Hessian

Let  $\mathbf{H}_0, \mathbf{H}_1, \mathbf{H}_2, \dots$  be successive approximations of the inverse  $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$  of the Hessian. We now derive a condition that the approximations should

satisfy, which forms the starting point for our subsequent discussion of quasi-Newton algorithms. To begin, suppose first that the Hessian matrix  $\mathbf{F}(\mathbf{x})$  of the objective function  $f$  is constant and independent of  $\mathbf{x}$ . In other words, the objective function is quadratic, with Hessian  $\mathbf{F}(\mathbf{x}) = \mathbf{Q}$  for all  $\mathbf{x}$ , where  $\mathbf{Q} = \mathbf{Q}^\top$ . Then,

$$\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)} = \mathbf{Q}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}).$$

Let

$$\Delta\mathbf{g}^{(k)} \triangleq \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}$$

and

$$\Delta\mathbf{x}^{(k)} \triangleq \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}.$$

Then, we may write

$$\Delta\mathbf{g}^{(k)} = \mathbf{Q}\Delta\mathbf{x}^{(k)}.$$

We start with a real symmetric positive definite matrix  $\mathbf{H}_0$ . Note that given  $k$ , the matrix  $\mathbf{Q}^{-1}$  satisfies

$$\mathbf{Q}^{-1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}, \quad 0 \leq i \leq k.$$

Therefore, we also impose the requirement that the approximation  $\mathbf{H}_{k+1}$  of the Hessian satisfy

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}, \quad 0 \leq i \leq k.$$

If  $n$  steps are involved, then moving in  $n$  directions  $\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}$  yields

$$\begin{aligned} \mathbf{H}_n\Delta\mathbf{g}^{(0)} &= \Delta\mathbf{x}^{(0)}, \\ \mathbf{H}_n\Delta\mathbf{g}^{(1)} &= \Delta\mathbf{x}^{(1)}, \\ &\vdots \\ \mathbf{H}_n\Delta\mathbf{g}^{(n-1)} &= \Delta\mathbf{x}^{(n-1)}. \end{aligned}$$

This set of equations can be represented as

$$\mathbf{H}_n[\Delta\mathbf{g}^{(0)}, \Delta\mathbf{g}^{(1)}, \dots, \Delta\mathbf{g}^{(n-1)}] = [\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}].$$

Note that  $\mathbf{Q}$  satisfies

$$\mathbf{Q}[\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}] = [\Delta\mathbf{g}^{(0)}, \Delta\mathbf{g}^{(1)}, \dots, \Delta\mathbf{g}^{(n-1)}]$$

and

$$\mathbf{Q}^{-1}[\Delta\mathbf{g}^{(0)}, \Delta\mathbf{g}^{(1)}, \dots, \Delta\mathbf{g}^{(n-1)}] = [\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}].$$

Therefore, if  $[\Delta\mathbf{g}^{(0)}, \Delta\mathbf{g}^{(1)}, \dots, \Delta\mathbf{g}^{(n-1)}]$  is nonsingular, then  $\mathbf{Q}^{-1}$  is determined uniquely after  $n$  steps, via

$$\mathbf{Q}^{-1} = \mathbf{H}_n = [\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}][\Delta\mathbf{g}^{(0)}, \Delta\mathbf{g}^{(1)}, \dots, \Delta\mathbf{g}^{(n-1)}]^{-1}.$$

As a consequence, we conclude that if  $\mathbf{H}_n$  satisfies the equations  $\mathbf{H}_n \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq n - 1$ , then the algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k \mathbf{g}^{(k)}$ ,  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)})$ , is guaranteed to solve problems with quadratic objective functions in  $n + 1$  steps, because the update  $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \alpha_n \mathbf{H}_n \mathbf{g}^{(n)}$  is equivalent to Newton's algorithm. In fact, as we shall see below (Theorem 11.1), such algorithms solve quadratic problems of  $n$  variables in at most  $n$  steps.

The considerations above illustrate the basic idea behind the quasi-Newton methods. Specifically, quasi-Newton algorithms have the form

$$\begin{aligned}\mathbf{d}^{(k)} &= -\mathbf{H}_k \mathbf{g}^{(k)}, \\ \alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}), \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},\end{aligned}$$

where the matrices  $\mathbf{H}_0, \mathbf{H}_1, \dots$  are symmetric. In the quadratic case these matrices are required to satisfy

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, \quad 0 \leq i \leq k,$$

where  $\Delta \mathbf{x}^{(i)} = \mathbf{x}^{(i+1)} - \mathbf{x}^{(i)} = \alpha_i \mathbf{d}^{(i)}$  and  $\Delta \mathbf{g}^{(i)} = \mathbf{g}^{(i+1)} - \mathbf{g}^{(i)} = \mathbf{Q} \Delta \mathbf{x}^{(i)}$ . It turns out that quasi-Newton methods are also conjugate direction methods, as stated in the following.

**Theorem 11.1** *Consider a quasi-Newton algorithm applied to a quadratic function with Hessian  $\mathbf{Q} = \mathbf{Q}^\top$  such that for  $0 \leq k < n - 1$ ,*

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, \quad 0 \leq i \leq k,$$

where  $\mathbf{H}_{k+1} = \mathbf{H}_{k+1}^\top$ . If  $\alpha_i \neq 0$ ,  $0 \leq i \leq k$ , then  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$  are  $\mathbf{Q}$ -conjugate.  $\square$

*Proof.* We proceed by induction. We begin with the  $k = 0$  case: that  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate. Because  $\alpha_0 \neq 0$ , we can write  $\mathbf{d}^{(0)} = \Delta \mathbf{x}^{(0)} / \alpha_0$ . Hence,

$$\begin{aligned}\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)} &= -\mathbf{g}^{(1)\top} \mathbf{H}_1 \mathbf{Q} \mathbf{d}^{(0)} \\ &= -\mathbf{g}^{(1)\top} \mathbf{H}_1 \frac{\mathbf{Q} \Delta \mathbf{x}^{(0)}}{\alpha_0} \\ &= -\mathbf{g}^{(1)\top} \frac{\mathbf{H}_1 \Delta \mathbf{g}^{(0)}}{\alpha_0} \\ &= -\mathbf{g}^{(1)\top} \frac{\Delta \mathbf{x}^{(0)}}{\alpha_0} \\ &= -\mathbf{g}^{(1)\top} \mathbf{d}^{(0)}.\end{aligned}$$

But  $\mathbf{g}^{(1)\top} \mathbf{d}^{(0)} = 0$  as a consequence of  $\alpha_0 > 0$  being the minimizer of  $\phi(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$  (see Exercise 11.1). Hence,  $\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)} = 0$ .

Assume that the result is true for  $k - 1$  (where  $k < n - 1$ ). We now prove the result for  $k$ , that is, that  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$  are  $\mathbf{Q}$ -conjugate. It suffices to show that  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ ,  $0 \leq i \leq k$ . Given  $i$ ,  $0 \leq i \leq k$ , using the same algebraic steps as in the  $k = 0$  case, and using the assumption that  $\alpha_i \neq 0$ , we obtain

$$\begin{aligned} \mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} &= -\mathbf{g}^{(k+1)\top} \mathbf{H}_{k+1} \mathbf{Q} \mathbf{d}^{(i)} \\ &\vdots \\ &= -\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)}. \end{aligned}$$

Because  $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$  are  $\mathbf{Q}$ -conjugate by assumption, we conclude from Lemma 10.2 that  $\mathbf{g}^{(k+1)\top} \mathbf{d}^{(i)} = 0$ . Hence,  $\mathbf{d}^{(k+1)\top} \mathbf{Q} \mathbf{d}^{(i)} = 0$ , which completes the proof. ■

By Theorem 11.1 we conclude that a quasi-Newton algorithm solves a quadratic of  $n$  variables in at most  $n$  steps.

Note that the equations that the matrices  $\mathbf{H}_k$  are required to satisfy do not determine those matrices uniquely. Thus, we have some freedom in the way we compute the  $\mathbf{H}_k$ . In the methods we describe, we compute  $\mathbf{H}_{k+1}$  by adding a correction to  $\mathbf{H}_k$ . In the following sections we consider three specific updating formulas.

### 11.3 The Rank One Correction Formula

In the *rank one correction formula*, the correction term is symmetric and has the form  $a_k \mathbf{z}^{(k)} \mathbf{z}^{(k)\top}$ , where  $a_k \in \mathbb{R}$  and  $\mathbf{z}^{(k)} \in \mathbb{R}^n$ . Therefore, the update equation is

$$\mathbf{H}_{k+1} = \mathbf{H}_k + a_k \mathbf{z}^{(k)} \mathbf{z}^{(k)\top}.$$

Note that

$$\text{rank } \mathbf{z}^{(k)} \mathbf{z}^{(k)\top} = \text{rank} \left( \begin{bmatrix} z_1^{(k)} \\ \vdots \\ z_n^{(k)} \end{bmatrix} \begin{bmatrix} z_1^{(k)} & \cdots & z_n^{(k)} \end{bmatrix} \right) = 1$$

and hence the name *rank one correction* [it is also called the *single-rank symmetric (SRS) algorithm*]. The product  $\mathbf{z}^{(k)} \mathbf{z}^{(k)\top}$  is sometimes referred to as the *dyadic product* or *outer product*. Observe that if  $\mathbf{H}_k$  is symmetric, then so is  $\mathbf{H}_{k+1}$ .

Our goal now is to determine  $a_k$  and  $\mathbf{z}^{(k)}$ , given  $\mathbf{H}_k$ ,  $\Delta \mathbf{g}^{(k)}$ ,  $\Delta \mathbf{x}^{(k)}$ , so that the required relationship discussed in Section 11.2 is satisfied; namely,

$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}$ ,  $i = 1, \dots, k$ . To begin, let us first consider the condition  $\mathbf{H}_{k+1}\Delta\mathbf{g}^{(k)} = \Delta\mathbf{x}^{(k)}$ . In other words, given  $\mathbf{H}_k$ ,  $\Delta\mathbf{g}^{(k)}$ , and  $\Delta\mathbf{x}^{(k)}$ , we wish to find  $a_k$  and  $\mathbf{z}^{(k)}$  to ensure that

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(k)} = (\mathbf{H}_k + a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)\top})\Delta\mathbf{g}^{(k)} = \Delta\mathbf{x}^{(k)}.$$

First note that  $\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)}$  is a scalar. Thus,

$$\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)} = (a_k\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})\mathbf{z}^{(k)},$$

and hence

$$\mathbf{z}^{(k)} = \frac{\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)}}{a_k(\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})}.$$

We can now determine

$$a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)\top} = \frac{(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})^\top}{a_k(\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})^2}.$$

Hence,

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})^\top}{a_k(\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})^2}.$$

The next step is to express the denominator of the second term on the right-hand side of the equation above as a function of the given quantities  $\mathbf{H}_k$ ,  $\Delta\mathbf{g}^{(k)}$ , and  $\Delta\mathbf{x}^{(k)}$ . To accomplish this, premultiply  $\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)} = (a_k\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})\mathbf{z}^{(k)}$  by  $\Delta\mathbf{g}^{(k)\top}$  to obtain

$$\Delta\mathbf{g}^{(k)\top}\Delta\mathbf{x}^{(k)} - \Delta\mathbf{g}^{(k)\top}\mathbf{H}_k\Delta\mathbf{g}^{(k)} = \Delta\mathbf{g}^{(k)\top}a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)}.$$

Observe that  $a_k$  is a scalar and so is  $\Delta\mathbf{g}^{(k)\top}\mathbf{z}^{(k)} = \mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)}$ . Thus,

$$\Delta\mathbf{g}^{(k)\top}\Delta\mathbf{x}^{(k)} - \Delta\mathbf{g}^{(k)\top}\mathbf{H}_k\Delta\mathbf{g}^{(k)} = a_k(\mathbf{z}^{(k)\top}\Delta\mathbf{g}^{(k)})^2.$$

Taking this relation into account yields

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})^\top}{\Delta\mathbf{g}^{(k)\top}(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})}.$$

We summarize the above development in the following algorithm.

### Rank One Algorithm

1. Set  $k := 0$ ; select  $\mathbf{x}^{(0)}$  and a real symmetric positive definite  $\mathbf{H}_0$ .
2. If  $\mathbf{g}^{(k)} = \mathbf{0}$ , stop; else,  $\mathbf{d}^{(k)} = -\mathbf{H}_k\mathbf{g}^{(k)}$ .

3. Compute

$$\begin{aligned}\alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}), \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}.\end{aligned}$$

4. Compute

$$\begin{aligned}\Delta \mathbf{x}^{(k)} &= \alpha_k \mathbf{d}^{(k)}, \\ \Delta \mathbf{g}^{(k)} &= \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}, \\ \mathbf{H}_{k+1} &= \mathbf{H}_k + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top}{\Delta \mathbf{g}^{(k)^\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}.\end{aligned}$$

5. Set  $k := k + 1$ ; go to step 2.

The rank one algorithm is based on satisfying the equation

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(k)} = \Delta \mathbf{x}^{(k)}.$$

However, what we want is

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, \quad i = 0, 1, \dots, k.$$

It turns out that the above is, in fact, true automatically, as stated in the following theorem.

**Theorem 11.2** *For the rank one algorithm applied to the quadratic with Hessian  $\mathbf{Q} = \mathbf{Q}^\top$ , we have  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ .*  $\square$

*Proof.* We prove the result by induction. From the discussion before the theorem, it is clear that the claim is true for  $k = 0$ . Suppose now that the theorem is true for  $k - 1 \geq 0$ ; that is,  $\mathbf{H}_k \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $i < k$ . We now show that the theorem is true for  $k$ . Our construction of the correction term ensures that

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(k)} = \Delta \mathbf{x}^{(k)}.$$

So we only have to show that

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, \quad i < k.$$

To this end, fix  $i < k$ . We have

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \mathbf{H}_k \Delta \mathbf{g}^{(i)} + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top}{\Delta \mathbf{g}^{(k)^\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})} \Delta \mathbf{g}^{(i)}.$$

By the induction hypothesis,  $\mathbf{H}_k \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ . To complete the proof, it is enough to show that the second term on the right-hand side of the equation above is equal to zero. For this to be true it is enough that

$$(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(k)^\top} \Delta \mathbf{g}^{(i)} - \Delta \mathbf{g}^{(k)^\top} \mathbf{H}_k \Delta \mathbf{g}^{(i)} = 0.$$

Indeed, since

$$\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(i)} = \Delta \mathbf{g}^{(k)\top} (\mathbf{H}_k \Delta \mathbf{g}^{(i)}) = \Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(i)}$$

by the induction hypothesis, and because  $\Delta \mathbf{g}^{(k)} = \mathbf{Q} \Delta \mathbf{x}^{(k)}$ , we have

$$\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(i)} = \Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(i)} = \Delta \mathbf{x}^{(k)\top} \mathbf{Q} \Delta \mathbf{x}^{(i)} = \Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)}.$$

Hence,

$$(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)} - \Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)} = 0,$$

which completes the proof. ■

**Example 11.1** Let

$$f(x_1, x_2) = x_1^2 + \frac{1}{2}x_2^2 + 3.$$

Apply the rank one correction algorithm to minimize  $f$ . Use  $\mathbf{x}^{(0)} = [1, 2]^\top$  and  $\mathbf{H}_0 = \mathbf{I}_2$  ( $2 \times 2$  identity matrix).

We can represent  $f$  as

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x} + 3.$$

Thus,

$$\mathbf{g}^{(k)} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x}^{(k)}.$$

Because  $\mathbf{H}_0 = \mathbf{I}_2$ ,

$$\mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = [-2, -2]^\top.$$

The objective function is quadratic, and hence

$$\begin{aligned} \alpha_0 &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}) = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} \\ &= \frac{[2, 2] \begin{bmatrix} 2 \\ 2 \end{bmatrix}}{[2, 2] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}} = \frac{2}{3}, \end{aligned}$$

and thus

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \left[ -\frac{1}{3}, \frac{2}{3} \right]^\top.$$

We then compute

$$\begin{aligned}\Delta \mathbf{x}^{(0)} &= \alpha_0 \mathbf{d}^{(0)} = \left[ -\frac{4}{3}, -\frac{4}{3} \right]^\top, \\ \mathbf{g}^{(1)} &= \mathbf{Q} \mathbf{x}^{(1)} = \left[ -\frac{2}{3}, \frac{2}{3} \right]^\top, \\ \Delta \mathbf{g}^{(0)} &= \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \left[ -\frac{8}{3}, -\frac{4}{3} \right]^\top.\end{aligned}$$

Because

$$\Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) = \left[ -\frac{8}{3}, -\frac{4}{3} \right] \left[ \frac{4}{3} \atop 0 \right] = -\frac{32}{9},$$

we obtain

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix}.$$

Therefore,

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \left[ \frac{1}{3}, -\frac{2}{3} \right]^\top$$

and

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 1.$$

We now compute

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [0, 0]^\top.$$

Note that  $\mathbf{g}^{(2)} = \mathbf{0}$ , and therefore  $\mathbf{x}^{(2)} = \mathbf{x}^*$ . As expected, the algorithm solves the problem in two steps.

Note that the directions  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate, in accordance with Theorem 11.1. ■

Unfortunately, the rank one correction algorithm is not very satisfactory, for several reasons. First, the matrix  $\mathbf{H}_{k+1}$  that the rank one algorithm generates may not be positive definite (see Example 11.2 below) and thus  $\mathbf{d}^{(k+1)}$  may not be a descent direction. This happens even in the quadratic case (see Example 11.10). Furthermore, if

$$\Delta \mathbf{g}^{(k)} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})$$

is close to zero, then there may be numerical problems in evaluating  $\mathbf{H}_{k+1}$ .

**Example 11.2** Assume that  $\mathbf{H}_k > 0$ . It turns out that if  $\Delta \mathbf{g}^{(k)\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) > 0$ , then  $\mathbf{H}_{k+1} > 0$  (see Exercise 11.7). However, if

$\Delta\mathbf{g}^{(k)\top}(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)}) < 0$ , then  $\mathbf{H}_{k+1}$  may not be positive definite. As an example of what might happen if  $\Delta\mathbf{g}^{(k)\top}(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)}) < 0$ , consider applying the rank one algorithm to the function

$$f(\mathbf{x}) = \frac{x_1^4}{4} + \frac{x_2^2}{2} - x_1x_2 + x_1 - x_2$$

with an initial point

$$\mathbf{x}^{(0)} = [0.59607, 0.59607]^\top$$

and initial matrix

$$\mathbf{H}_0 = \begin{bmatrix} 0.94913 & 0.14318 \\ 0.14318 & 0.59702 \end{bmatrix}.$$

Note that  $\mathbf{H}_0 > 0$ . We have

$$\Delta\mathbf{g}^{(0)\top}(\Delta\mathbf{x}^{(0)} - \mathbf{H}_0\Delta\mathbf{g}^{(0)}) = -0.03276$$

and

$$\mathbf{H}_1 = \begin{bmatrix} 0.94481 & 0.23324 \\ 0.23324 & -1.2788 \end{bmatrix}.$$

It is easy to check that  $\mathbf{H}_1$  is not positive definite (it is indefinite, with eigenvalues 0.96901 and  $-1.3030$ ). ■

Fortunately, alternative algorithms have been developed for updating  $\mathbf{H}_k$ . In particular, if we use a “rank two” update, then  $\mathbf{H}_k$  is guaranteed to be positive definite for all  $k$ , provided that the line search is exact. We discuss this in the next section.

## 11.4 The DFP Algorithm

The rank two update was originally developed by Davidon in 1959 and was subsequently modified by Fletcher and Powell in 1963: hence the name *DFP algorithm*. The DFP algorithm is also known as the *variable metric algorithm*. We summarize the algorithm below.

### DFP Algorithm

1. Set  $k := 0$ ; select  $\mathbf{x}^{(0)}$  and a real symmetric positive definite  $\mathbf{H}_0$ .
2. If  $\mathbf{g}^{(k)} = \mathbf{0}$ , stop; else,  $\mathbf{d}^{(k)} = -\mathbf{H}_k\mathbf{g}^{(k)}$ .
3. Compute

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha\mathbf{d}^{(k)}),$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}.$$

4. Compute

$$\begin{aligned}\Delta \mathbf{x}^{(k)} &= \alpha_k \mathbf{d}^{(k)}, \\ \Delta \mathbf{g}^{(k)} &= \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}, \\ \mathbf{H}_{k+1} &= \mathbf{H}_k + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{[\mathbf{H}_k \Delta \mathbf{g}^{(k)}][\mathbf{H}_k \Delta \mathbf{g}^{(k)}]^\top}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}.\end{aligned}$$

5. Set  $k := k + 1$ ; go to step 2.

We now show that the DFP algorithm is a quasi-Newton method, in the sense that when applied to quadratic problems, we have  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ .

**Theorem 11.3** *In the DFP algorithm applied to the quadratic with Hessian  $\mathbf{Q} = \mathbf{Q}^\top$ , we have  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ .  $\square$*

*Proof.* We use induction. For  $k = 0$ , we have

$$\begin{aligned}\mathbf{H}_1 \Delta \mathbf{g}^{(0)} &= \mathbf{H}_0 \Delta \mathbf{g}^{(0)} + \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top}}{\Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)}} \Delta \mathbf{g}^{(0)} - \frac{\mathbf{H}_0 \Delta \mathbf{g}^{(0)} \Delta \mathbf{g}^{(0)\top} \mathbf{H}_0}{\Delta \mathbf{g}^{(0)\top} \mathbf{H}_0 \Delta \mathbf{g}^{(0)}} \Delta \mathbf{g}^{(0)} \\ &= \Delta \mathbf{x}^{(0)}.\end{aligned}$$

Assume that the result is true for  $k - 1$ ; that is,  $\mathbf{H}_k \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k - 1$ . We now show that  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . First, consider  $i = k$ . We have

$$\begin{aligned}\mathbf{H}_{k+1} \Delta \mathbf{g}^{(k)} &= \mathbf{H}_k \Delta \mathbf{g}^{(k)} + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} \Delta \mathbf{g}^{(k)} - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \Delta \mathbf{g}^{(k)} \\ &= \Delta \mathbf{x}^{(k)}.\end{aligned}$$

It remains to consider the case  $i < k$ . To this end,

$$\begin{aligned}\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} &= \mathbf{H}_k \Delta \mathbf{g}^{(i)} + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} \Delta \mathbf{g}^{(i)} - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \Delta \mathbf{g}^{(i)} \\ &= \Delta \mathbf{x}^{(i)} + \frac{\Delta \mathbf{x}^{(k)}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} (\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)}) \\ &\quad - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} (\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(i)}).\end{aligned}$$

Now,

$$\begin{aligned}\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(i)} &= \Delta \mathbf{x}^{(k)\top} \mathbf{Q} \Delta \mathbf{x}^{(i)} \\ &= \alpha_k \alpha_i \mathbf{d}^{(k)\top} \mathbf{Q} \mathbf{d}^{(i)} \\ &= 0,\end{aligned}$$

by the induction hypothesis and Theorem 11.1. The same arguments yield  $\Delta\mathbf{g}^{(k)\top}\Delta\mathbf{x}^{(i)} = 0$ . Hence,

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)},$$

which completes the proof.  $\blacksquare$

By Theorems 11.1 and 11.3 we conclude that the DFP algorithm is a conjugate direction algorithm.

**Example 11.3** Locate the minimizer of

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \quad \mathbf{x} \in \mathbb{R}^2.$$

Use the initial point  $\mathbf{x}^{(0)} = [0, 0]^\top$  and  $\mathbf{H}_0 = \mathbf{I}_2$ .

Note that in this case

$$\mathbf{g}^{(k)} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x}^{(k)} - \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

Hence,

$$\begin{aligned} \mathbf{g}^{(0)} &= [1, -1]^\top, \\ \mathbf{d}^{(0)} &= -\mathbf{H}_0\mathbf{g}^{(0)} = -\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}. \end{aligned}$$

Because  $f$  is a quadratic function,

$$\begin{aligned} \alpha_0 &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha\mathbf{d}^{(0)}) = -\frac{\mathbf{g}^{(0)\top}\mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top}\mathbf{Q}\mathbf{d}^{(0)}} \\ &= -\frac{[1, -1] \begin{bmatrix} -1 \\ 1 \end{bmatrix}}{[-1, 1] \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix}} = 1. \end{aligned}$$

Therefore,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0\mathbf{d}^{(0)} = [-1, 1]^\top.$$

We then compute

$$\begin{aligned} \Delta\mathbf{x}^{(0)} &= \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = [-1, 1]^\top, \\ \mathbf{g}^{(1)} &= \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \end{aligned}$$

and

$$\Delta \mathbf{g}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [-2, 0]^\top.$$

Observe that

$$\begin{aligned}\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top} &= \begin{bmatrix} -1 \\ 1 \end{bmatrix} [-1, 1] = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}, \\ \Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)} &= [-1, 1] \begin{bmatrix} -2 \\ 0 \end{bmatrix} = 2, \\ \mathbf{H}_0 \Delta \mathbf{g}^{(0)} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 0 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \end{bmatrix}.\end{aligned}$$

Thus,

$$(\mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top = \begin{bmatrix} -2 \\ 0 \end{bmatrix} [-2, 0] = \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix}$$

and

$$\Delta \mathbf{g}^{(0)\top} \mathbf{H}_0 \Delta \mathbf{g}^{(0)} = [-2, 0] \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 0 \end{bmatrix} = 4.$$

Using the above, we now compute  $\mathbf{H}_1$ :

$$\begin{aligned}\mathbf{H}_1 &= \mathbf{H}_0 + \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top}}{\Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)}} - \frac{(\mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\mathbf{H}_0 \Delta \mathbf{g}^{(0)})^\top}{\Delta \mathbf{g}^{(0)\top} \mathbf{H}_0 \Delta \mathbf{g}^{(0)}} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} - \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{3}{2} \end{bmatrix}.\end{aligned}$$

We now compute  $\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = [0, 1]^\top$  and

$$\alpha_1 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(1)} + \alpha \mathbf{d}^{(1)}) = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = \frac{1}{2}.$$

Hence,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [-1, 3/2]^\top = \mathbf{x}^*,$$

because  $f$  is a quadratic function of two variables.

Note that we have  $\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(1)} = \mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)} = 0$ ; that is,  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate directions. ■

We now show that in the DFP algorithm,  $\mathbf{H}_{k+1}$  inherits positive definiteness from  $\mathbf{H}_k$ .

**Theorem 11.4** Suppose that  $\mathbf{g}^{(k)} \neq \mathbf{0}$ . In the DFP algorithm, if  $\mathbf{H}_k$  is positive definite, then so is  $\mathbf{H}_{k+1}$ .  $\square$

*Proof.* We first write the following quadratic form:

$$\begin{aligned}\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} &= \mathbf{x}^\top \mathbf{H}_k \mathbf{x} + \frac{\mathbf{x}^\top \Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top} \mathbf{x}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{\mathbf{x}^\top (\mathbf{H}_k \Delta \mathbf{g}^{(k)}) (\mathbf{H}_k \Delta \mathbf{g}^{(k)})^\top \mathbf{x}}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}} \\ &= \mathbf{x}^\top \mathbf{H}_k \mathbf{x} + \frac{(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{(\mathbf{x}^\top \mathbf{H}_k \Delta \mathbf{g}^{(k)})^2}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}.\end{aligned}$$

Define

$$\begin{aligned}\mathbf{a} &\triangleq \mathbf{H}_k^{1/2} \mathbf{x}, \\ \mathbf{b} &\triangleq \mathbf{H}_k^{1/2} \Delta \mathbf{g}^{(k)},\end{aligned}$$

where

$$\mathbf{H}_k = \mathbf{H}_k^{1/2} \mathbf{H}_k^{1/2}.$$

Note that because  $\mathbf{H}_k > 0$ , its square root is well-defined; see Section 3.4 for more information on this property of positive definite matrices. Using the definitions of  $\mathbf{a}$  and  $\mathbf{b}$ , we obtain

$$\begin{aligned}\mathbf{x}^\top \mathbf{H}_k \mathbf{x} &= \mathbf{x}^\top \mathbf{H}_k^{1/2} \mathbf{H}_k^{1/2} \mathbf{x} = \mathbf{a}^\top \mathbf{a}, \\ \mathbf{x}^\top \mathbf{H}_k \Delta \mathbf{g}^{(k)} &= \mathbf{x}^\top \mathbf{H}_k^{1/2} \mathbf{H}_k^{1/2} \Delta \mathbf{g}^{(k)} = \mathbf{a}^\top \mathbf{b},\end{aligned}$$

and

$$\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)} = \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k^{1/2} \mathbf{H}_k^{1/2} \Delta \mathbf{g}^{(k)} = \mathbf{b}^\top \mathbf{b}.$$

Hence,

$$\begin{aligned}\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} &= \mathbf{a}^\top \mathbf{a} + \frac{(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{(\mathbf{a}^\top \mathbf{b})^2}{\mathbf{b}^\top \mathbf{b}} \\ &= \frac{\|\mathbf{a}\|^2 \|\mathbf{b}\|^2 - (\langle \mathbf{a}, \mathbf{b} \rangle)^2}{\|\mathbf{b}\|^2} + \frac{(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}}.\end{aligned}$$

We also have

$$\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)} = \Delta \mathbf{x}^{(k)\top} (\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}) = -\Delta \mathbf{x}^{(k)\top} \mathbf{g}^{(k)},$$

since  $\Delta \mathbf{x}^{(k)\top} \mathbf{g}^{(k+1)} = \alpha_k \mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} = 0$  by Lemma 10.2 (see also Exercise 11.1). Because

$$\Delta \mathbf{x}^{(k)} = \alpha_k \mathbf{d}^{(k)} = -\alpha_k \mathbf{H}_k \mathbf{g}^{(k)},$$

we have

$$\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)} = -\Delta \mathbf{x}^{(k)\top} \mathbf{g}^{(k)} = \alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}.$$

This yields

$$\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} = \frac{\|\mathbf{a}\|^2 \|\mathbf{b}\|^2 - (\langle \mathbf{a}, \mathbf{b} \rangle)^2}{\|\mathbf{b}\|^2} + \frac{(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2}{\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}}.$$

Both terms on the right-hand side of the above equation are nonnegative—the first term is nonnegative because of the Cauchy-Schwarz inequality, and the second term is nonnegative because  $\mathbf{H}_k > 0$  and  $\alpha_k > 0$  (by Proposition 11.1). Therefore, to show that  $\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} > 0$  for  $\mathbf{x} \neq \mathbf{0}$ , we only need to demonstrate that these terms do not both vanish simultaneously.

The first term vanishes only if  $\mathbf{a}$  and  $\mathbf{b}$  are proportional, that is, if  $\mathbf{a} = \beta \mathbf{b}$  for some scalar  $\beta$ . Thus, to complete the proof it is enough to show that if  $\mathbf{a} = \beta \mathbf{b}$ , then  $(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2 / (\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}) > 0$ . Indeed, first observe that

$$\mathbf{H}_k^{1/2} \mathbf{x} = \mathbf{a} = \beta \mathbf{b} = \beta \mathbf{H}_k^{1/2} \Delta \mathbf{g}^{(k)} = \mathbf{H}_k^{1/2} (\beta \Delta \mathbf{g}^{(k)}).$$

Hence,

$$\mathbf{x} = \beta \Delta \mathbf{g}^{(k)}.$$

Using the expression for  $\mathbf{x}$  above and the expression  $\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)} = \alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}$ , we obtain

$$\begin{aligned} \frac{(\mathbf{x}^\top \Delta \mathbf{x}^{(k)})^2}{\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}} &= \frac{\beta^2 (\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)})^2}{\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}} = \frac{\beta^2 (\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)})^2}{\alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)}} \\ &= \beta^2 \alpha_k \mathbf{g}^{(k)\top} \mathbf{H}_k \mathbf{g}^{(k)} > 0. \end{aligned}$$

Thus, for all  $\mathbf{x} \neq \mathbf{0}$ ,

$$\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} > 0,$$

which completes the proof. ■

The DFP algorithm is superior to the rank one algorithm in that it preserves the positive definiteness of  $\mathbf{H}_k$ . However, it turns out that in the case of larger nonquadratic problems the algorithm has the tendency of sometimes getting “stuck.” This phenomenon is attributed to  $\mathbf{H}_k$  becoming nearly singular [19]. In the next section we discuss an algorithm that alleviates this problem.

## 11.5 The BFGS Algorithm

In 1970, an alternative update formula was suggested independently by Broyden, Fletcher, Goldfarb, and Shanno. The method, now called the *BFGS algorithm*, is discussed in this section.

To derive the BFGS update, we use the concept of *duality*, or *complementarity*, as presented in [43] and [88]. To discuss this concept, recall that the

updating formulas for the approximation of the inverse of the Hessian matrix were based on satisfying the equations

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, \quad 0 \leq i \leq k,$$

which were derived from  $\Delta \mathbf{g}^{(i)} = \mathbf{Q} \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . We then formulated update formulas for the approximations to the inverse of the Hessian matrix  $\mathbf{Q}^{-1}$ . An alternative to approximating  $\mathbf{Q}^{-1}$  is to approximate  $\mathbf{Q}$  itself. To do this let  $\mathbf{B}_k$  be our estimate of  $\mathbf{Q}$  at the  $k$ th step. We require  $\mathbf{B}_{k+1}$  to satisfy

$$\Delta \mathbf{g}^{(i)} = \mathbf{B}_{k+1} \Delta \mathbf{x}^{(i)}, \quad 0 \leq i \leq k.$$

Notice that this set of equations is similar to the previous set of equations for  $\mathbf{H}_{k+1}$ , the only difference being that the roles of  $\Delta \mathbf{x}^{(i)}$  and  $\Delta \mathbf{g}^{(i)}$  are interchanged. Thus, given any update formula for  $\mathbf{H}_k$ , a corresponding update formula for  $\mathbf{B}_k$  can be found by interchanging the roles of  $\mathbf{B}_k$  and  $\mathbf{H}_k$  and of  $\Delta \mathbf{g}^{(k)}$  and  $\Delta \mathbf{x}^{(k)}$ . In particular, the BFGS update for  $\mathbf{B}_k$  corresponds to the DFP update for  $\mathbf{H}_k$ . Formulas related in this way are said to be *dual* or *complementary* [43].

Recall that the DFP update for the approximation  $\mathbf{H}_k$  of the inverse Hessian is

$$\mathbf{H}_{k+1}^{DFP} = \mathbf{H}_k + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top}}{\Delta \mathbf{x}^{(k)\top} \Delta \mathbf{g}^{(k)}} - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k}{\Delta \mathbf{g}^{(k)\top} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}.$$

Using the complementarity concept, we can easily obtain an update equation for the approximation  $\mathbf{B}_k$  of the Hessian:

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}} - \frac{\mathbf{B}_k \Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)}}.$$

This is the BFGS update of  $\mathbf{B}_k$ .

Now, to obtain the BFGS update for the approximation of the inverse Hessian, we take the inverse of  $\mathbf{B}_{k+1}$  to obtain

$$\begin{aligned} \mathbf{H}_{k+1}^{BFGS} &= (\mathbf{B}_{k+1})^{-1} \\ &= \left( \mathbf{B}_k + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)\top}}{\Delta \mathbf{g}^{(k)\top} \Delta \mathbf{x}^{(k)}} - \frac{\mathbf{B}_k \Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)\top} \mathbf{B}_k}{\Delta \mathbf{x}^{(k)\top} \mathbf{B}_k \Delta \mathbf{x}^{(k)}} \right)^{-1}. \end{aligned}$$

To compute  $\mathbf{H}_{k+1}^{BFGS}$  by inverting the right-hand side of this equation, we apply the following formula for a matrix inverse, known as the *Sherman-Morrison formula* (see [63, p. 123] or [53, p. 50]).

**Lemma 11.1** *Let  $\mathbf{A}$  be a nonsingular matrix. Let  $\mathbf{u}$  and  $\mathbf{v}$  be column vectors such that  $1 + \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u} \neq 0$ . Then,  $\mathbf{A} + \mathbf{u}\mathbf{v}^\top$  is nonsingular, and its inverse can be written in terms of  $\mathbf{A}^{-1}$  using the following formula:*

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{(\mathbf{A}^{-1}\mathbf{u})(\mathbf{v}^\top \mathbf{A}^{-1})}{1 + \mathbf{v}^\top \mathbf{A}^{-1}\mathbf{u}}.$$

□

*Proof.* We can prove the result easily by verification. ■

From Lemma 11.1 it follows that if  $\mathbf{A}^{-1}$  is known, then the inverse of the matrix  $\mathbf{A}$  augmented by a rank one matrix can be obtained by a modification of the matrix  $\mathbf{A}^{-1}$ .

Applying Lemma 11.1 twice to  $\mathbf{B}_{k+1}$  (see Exercise 11.12) yields

$$\begin{aligned} \mathbf{H}_{k+1}^{BFGS} &= \mathbf{H}_k + \left(1 + \frac{\Delta\mathbf{g}^{(k)\top} \mathbf{H}_k \Delta\mathbf{g}^{(k)}}{\Delta\mathbf{g}^{(k)\top} \Delta\mathbf{x}^{(k)}}\right) \frac{\Delta\mathbf{x}^{(k)} \Delta\mathbf{x}^{(k)\top}}{\Delta\mathbf{x}^{(k)\top} \Delta\mathbf{g}^{(k)}} \\ &\quad - \frac{\mathbf{H}_k \Delta\mathbf{g}^{(k)} \Delta\mathbf{x}^{(k)\top} + (\mathbf{H}_k \Delta\mathbf{g}^{(k)} \Delta\mathbf{x}^{(k)\top})^\top}{\Delta\mathbf{g}^{(k)\top} \Delta\mathbf{x}^{(k)}}, \end{aligned}$$

which represents the BFGS formula for updating  $\mathbf{H}_k$ .

Recall that for the quadratic case the DFP algorithm satisfies  $\mathbf{H}_{k+1}^{DFP} \Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . Therefore, the BFGS update for  $\mathbf{B}_k$  satisfies  $\mathbf{B}_{k+1} \Delta\mathbf{x}^{(i)} = \Delta\mathbf{g}^{(i)}$ ,  $0 \leq i \leq k$ . By construction of the BFGS formula for  $\mathbf{H}_{k+1}^{BFGS}$ , we conclude that  $\mathbf{H}_{k+1}^{BFGS} \Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ . Hence, the BFGS algorithm enjoys all the properties of quasi-Newton methods, including the conjugate directions property. Moreover, the BFGS algorithm also inherits the positive definiteness property of the DFP algorithm; that is, if  $\mathbf{g}^{(k)} \neq \mathbf{0}$  and  $\mathbf{H}_k > 0$ , then  $\mathbf{H}_{k+1}^{BFGS} > 0$ .

The BFGS update is reasonably robust when the line searches are sloppy (see [19]). This property allows us to save time in the line search part of the algorithm. The BFGS formula is often far more efficient than the DFP formula (see [107] for further discussion).

We conclude our discussion of the BFGS algorithm with the following numerical example.

**Example 11.4** Use the BFGS method to minimize

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b} + \log(\pi),$$

where

$$\mathbf{Q} = \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

Take  $\mathbf{H}_0 = \mathbf{I}_2$  and  $\mathbf{x}^{(0)} = [0, 0]^\top$ . Verify that  $\mathbf{H}_2 = \mathbf{Q}^{-1}$ .

We have

$$\mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = -(\mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}) = \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

The objective function is a quadratic, and hence we can use the following formula to compute  $\alpha_0$ :

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} Q \mathbf{d}^{(0)}} = \frac{1}{2}.$$

Therefore,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}.$$

To compute  $\mathbf{H}_1 = \mathbf{H}_1^{BFGS}$ , we need the following quantities:

$$\begin{aligned} \Delta \mathbf{x}^{(0)} &= \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}, \\ \mathbf{g}^{(1)} &= Q \mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} -3/2 \\ 0 \end{bmatrix}, \\ \Delta \mathbf{g}^{(0)} &= \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \begin{bmatrix} -3/2 \\ 1 \end{bmatrix}. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbf{H}_1 &= \mathbf{H}_0 + \left( 1 + \frac{\Delta \mathbf{g}^{(0)\top} \mathbf{H}_0 \Delta \mathbf{g}^{(0)}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \right) \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)\top}}{\Delta \mathbf{x}^{(0)\top} \Delta \mathbf{g}^{(0)}} \\ &\quad - \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{g}^{(0)\top} \mathbf{H}_0 + \mathbf{H}_0 \Delta \mathbf{g}^{(0)} \Delta \mathbf{x}^{(0)\top}}{\Delta \mathbf{g}^{(0)\top} \Delta \mathbf{x}^{(0)}} \\ &= \begin{bmatrix} 1 & 3/2 \\ 3/2 & 11/4 \end{bmatrix}. \end{aligned}$$

Hence, we have

$$\begin{aligned} \mathbf{d}^{(1)} &= -\mathbf{H}_1 \mathbf{g}^{(1)} = \begin{bmatrix} 3/2 \\ 9/4 \end{bmatrix}, \\ \alpha_1 &= -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} Q \mathbf{d}^{(1)}} = 2. \end{aligned}$$

Therefore,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}.$$

Because our objective function is a quadratic on  $\mathbb{R}^2$ ,  $\mathbf{x}^{(2)}$  is the minimizer. Notice that the gradient at  $\mathbf{x}^{(2)}$  is  $\mathbf{0}$ ; that is,  $\mathbf{g}^{(2)} = \mathbf{0}$ .

To verify that  $\mathbf{H}_2 = \mathbf{Q}^{-1}$ , we compute

$$\begin{aligned}\Delta\mathbf{x}^{(1)} &= \mathbf{x}^{(2)} - \mathbf{x}^{(1)} = \begin{bmatrix} 3 \\ 9/2 \end{bmatrix}, \\ \Delta\mathbf{g}^{(1)} &= \mathbf{g}^{(2)} - \mathbf{g}^{(1)} = \begin{bmatrix} 3/2 \\ 0 \end{bmatrix}.\end{aligned}$$

Hence,

$$\begin{aligned}\mathbf{H}_2 &= \mathbf{H}_1 + \left(1 + \frac{\Delta\mathbf{g}^{(1)\top} \mathbf{H}_1 \Delta\mathbf{g}^{(1)}}{\Delta\mathbf{g}^{(1)\top} \Delta\mathbf{x}^{(1)}}\right) \frac{\Delta\mathbf{x}^{(1)} \Delta\mathbf{x}^{(1)\top}}{\Delta\mathbf{x}^{(1)\top} \Delta\mathbf{g}^{(1)}} \\ &\quad - \frac{\Delta\mathbf{x}^{(1)} \Delta\mathbf{g}^{(1)\top} \mathbf{H}_1 + \mathbf{H}_1 \Delta\mathbf{g}^{(1)} \Delta\mathbf{x}^{(1)\top}}{\Delta\mathbf{g}^{(1)\top} \Delta\mathbf{x}^{(1)}} \\ &= \begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix}.\end{aligned}$$

Note that indeed  $\mathbf{H}_2 \mathbf{Q} = \mathbf{Q} \mathbf{H}_2 = \mathbf{I}_2$ , and hence  $\mathbf{H}_2 = \mathbf{Q}^{-1}$ . ■

For nonquadratic problems, quasi-Newton algorithms will not usually converge in  $n$  steps. As in the case of the conjugate gradient methods, here, too, some modifications may be necessary to deal with nonquadratic problems. For example, we may reinitialize the direction vector to the negative gradient after every few iterations (e.g.,  $n$  or  $n+1$ ), and continue until the algorithm satisfies the stopping criterion.

## EXERCISES

**11.1** Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)},$$

where  $\mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots$  are vectors in  $\mathbb{R}^n$ , and  $\alpha_k \geq 0$  is chosen to minimize  $f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ ; that is,

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

Note that the general algorithm above encompasses almost all algorithms that we discussed in this part, including the steepest descent, Newton, conjugate gradient, and quasi-Newton algorithms.

Let  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ , and assume that  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} < 0$ .

- a. Show that  $\mathbf{d}^{(k)}$  is a descent direction for  $f$  in the sense that there exists  $\bar{\alpha} > 0$  such that for all  $\alpha \in (0, \bar{\alpha}]$ ,

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}).$$

- b. Show that  $\alpha_k > 0$ .
- c. Show that  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k+1)} = 0$ .
- d. Show that the following algorithms all satisfy the condition  $\mathbf{d}^{(k)\top} \mathbf{g}^{(k)} < 0$ , if  $\mathbf{g}^{(k)} \neq \mathbf{0}$ :
1. Steepest descent algorithm.
  2. Newton's method, assuming that the Hessian is positive definite.
  3. Conjugate gradient algorithm.
  4. Quasi-Newton algorithm, assuming that  $\mathbf{H}_k > 0$ .
- e. For the case where  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b}$ , with  $\mathbf{Q} = \mathbf{Q}^\top > 0$ , derive an expression for  $\alpha_k$  in terms of  $\mathbf{Q}$ ,  $\mathbf{d}^{(k)}$ , and  $\mathbf{g}^{(k)}$ .

**11.2** Consider Newton's algorithm applied to a function  $f \in \mathcal{C}^2$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)}),$$

where  $\alpha_k$  is chosen according to a line search. Is this algorithm a member of the quasi-Newton family?

**11.3** In some optimization methods, when minimizing a given function  $f(\mathbf{x})$ , we select an initial guess  $\mathbf{x}^{(0)}$  and a real symmetric positive definite matrix  $\mathbf{H}_0$ . Then we iteratively compute  $\mathbf{H}_k$ ,  $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$  (where  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ ), and  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ , where

$$\alpha_k = \arg \min_{\alpha \geq 0} f\left(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}\right).$$

Suppose that the function we wish to minimize is a standard quadratic of the form

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b} + c, \quad \mathbf{Q} = \mathbf{Q}^\top > 0.$$

- a. Find an expression for  $\alpha_k$  in terms of  $\mathbf{Q}$ ,  $\mathbf{H}_k$ ,  $\mathbf{g}^{(k)}$ , and  $\mathbf{d}^{(k)}$ ;
- b. Give a sufficient condition on  $\mathbf{H}_k$  for  $\alpha_k$  to be positive.

**11.4** Consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H} \mathbf{g}^{(k)},$$

where, as usual,  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$  and  $\mathbf{H}$  is a fixed symmetric matrix.

- a. Suppose that  $f \in \mathcal{C}^3$  and there is a point  $\mathbf{x}^*$  such that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  and  $\mathbf{F}(\mathbf{x}^*)^{-1}$  exists. Find  $\mathbf{H}$  such that if  $\mathbf{x}^{(0)}$  is sufficiently close to  $\mathbf{x}^*$ , then  $\mathbf{x}^{(k)}$  converges to  $\mathbf{x}^*$  with order of convergence of at least 2.
- b. With the setting of  $\mathbf{H}$  in part a, is the given algorithm a quasi-Newton method?

### 11.5 Minimize the function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 1 \\ -1 \end{bmatrix} + 7$$

using the rank one correction method with the starting point  $\mathbf{x}^{(0)} = \mathbf{0}$ .

### 11.6 Consider the algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{M}_k \nabla f(\mathbf{x}^{(k)}),$$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ ,  $\mathbf{M}_k \in \mathbb{R}^{2 \times 2}$  is given by

$$\mathbf{M}_k = \begin{bmatrix} 1 & 0 \\ 0 & a \end{bmatrix}$$

with  $a \in \mathbb{R}$ , and

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{M}_k \nabla f(\mathbf{x}^{(k)})).$$

Suppose that at some iteration  $k$  we have  $\nabla f(\mathbf{x}^{(k)}) = [1, 1]^\top$ . Find the largest range of values of  $a$  that guarantees that  $\alpha_k > 0$  for any  $f$ .

**11.7** Consider the rank one algorithm. Assume that  $\mathbf{H}_k > 0$ . Show that if  $\Delta \mathbf{g}^{(k)\top} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}) > 0$ , then  $\mathbf{H}_{k+1} > 0$ .

**11.8** Based on the rank one update equation, derive an update formula using complementarity and the matrix inverse formula.

**11.9** Let

$$\begin{aligned} f &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b} + c \\ &= \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 1 \\ -1 \end{bmatrix} + 7 \end{aligned}$$

and  $\mathbf{x}^{(0)} = \mathbf{0}$ . Use the rank one correction method to generate two  $\mathbf{Q}$ -conjugate directions.

**11.10** Apply the rank one algorithm to the problem in Example 11.3.

**11.11** Consider the DFP algorithm applied to the quadratic function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ .

- a. Write down a formula for  $\alpha_k$  in terms of  $\mathbf{Q}$ ,  $\mathbf{g}^{(k)}$ , and  $\mathbf{d}^{(k)}$ .
- b. Show that if  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , then  $\alpha_k > 0$ .

**11.12** Use Lemma 11.1 to derive the BFGS update formula based on the DFP formula, using complementarity.

*Hint:* Define

$$\begin{aligned} \mathbf{A}_0 &= \mathbf{B}_k, \\ \mathbf{u}_0 &= \frac{\Delta\mathbf{g}^{(k)}}{\Delta\mathbf{g}^{(k)\top}\Delta\mathbf{x}^{(k)}}, \\ \mathbf{v}_0^\top &= \Delta\mathbf{g}^{(k)\top}, \\ \mathbf{u}_1 &= -\frac{\mathbf{B}_k\Delta\mathbf{x}^{(k)}}{\Delta\mathbf{x}^{(k)\top}\mathbf{B}_k\Delta\mathbf{x}^{(k)}}, \\ \mathbf{v}_1^\top &= \Delta\mathbf{x}^{(k)\top}\mathbf{B}_k, \\ \mathbf{A}_1 &= \mathbf{B}_k + \frac{\Delta\mathbf{g}^{(k)}\Delta\mathbf{g}^{(k)\top}}{\Delta\mathbf{g}^{(k)\top}\Delta\mathbf{x}^{(k)}} = \mathbf{A}_0 + \mathbf{u}_0\mathbf{v}_0^\top. \end{aligned}$$

Using the notation above, represent  $\mathbf{B}_{k+1}$  as

$$\begin{aligned} \mathbf{B}_{k+1} &= \mathbf{A}_0 + \mathbf{u}_0\mathbf{v}_0^\top + \mathbf{u}_1\mathbf{v}_1^\top \\ &= \mathbf{A}_1 + \mathbf{u}_1\mathbf{v}_1^\top. \end{aligned}$$

Apply Lemma 11.1 to the above.

**11.13** Assuming exact line search, show that if  $\mathbf{H}_0 = \mathbf{I}_n$  ( $n \times n$  identity matrix), then the first two steps of the BFGS algorithm yield the same points  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  as conjugate gradient algorithms with the Hestenes-Stiefel, the Polak-Ribi  re, and the Fletcher-Reeves formulas.

**11.14** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be such that  $f \in \mathcal{C}^1$ . Consider an optimization algorithm applied to this  $f$ , of the usual form  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ , where

$\alpha_k \geq 0$  is chosen according to line search. Suppose that  $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$ , where  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$  and  $\mathbf{H}_k$  is symmetric.

- a. Show that if  $\mathbf{H}_k$  satisfies the following conditions whenever the algorithm is applied to a quadratic, then the algorithm is quasi-Newton:
  1.  $\mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{U}_k$ .
  2.  $\mathbf{U}_k \Delta \mathbf{g}^{(k)} = \Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)}$ .
  3.  $\mathbf{U}_k = \mathbf{a}^{(k)} \Delta \mathbf{x}^{(k)\top} + \mathbf{b}^{(k)} \Delta \mathbf{g}^{(k)\top} \mathbf{H}_k$ , where  $\mathbf{a}^{(k)}$  and  $\mathbf{b}^{(k)}$  are in  $\mathbb{R}^n$ .
- b. Which (if any) among the rank-one, DFP, and BFGS algorithms satisfy the three conditions in part a (whenever the algorithm is applied to a quadratic)? For those that do, specify the vectors  $\mathbf{a}^{(k)}$  and  $\mathbf{b}^{(k)}$ .

**11.15** Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , consider an algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k \mathbf{g}^{(k)}$  for finding the minimizer of  $f$ , where  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$  and  $\mathbf{H}_k \in \mathbb{R}^{n \times n}$  is symmetric. Suppose that  $\mathbf{H}_k = \phi \mathbf{H}_k^{DFP} + (1 - \phi) \mathbf{H}_k^{BFGS}$ , where  $\phi \in \mathbb{R}$ , and  $\mathbf{H}_k^{DFP}$  and  $\mathbf{H}_k^{BFGS}$  are matrices generated by the DFP and BFGS algorithms, respectively.

- a. Show that the algorithm above is a quasi-Newton algorithm. Is the above algorithm a conjugate direction algorithm?
- b. Suppose that  $0 \leq \phi \leq 1$ . Show that if  $\mathbf{H}_0^{DFP} > 0$  and  $\mathbf{H}_0^{BFGS} > 0$ , then  $\mathbf{H}_k > 0$  for all  $k$ . What can you conclude from this about whether or not the algorithm has the descent property?

**11.16** Consider the following simple modification of the quasi-Newton family of algorithms. In the quadratic case, instead of the usual quasi-Newton condition  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ , suppose that we have  $\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \rho_i \Delta \mathbf{x}^{(i)}$ ,  $0 \leq i \leq k$ , where  $\rho_i > 0$ . We refer to the set of algorithms that satisfy the condition above as the *symmetric Huang family*.

Show that the symmetric Huang family algorithms are conjugate direction algorithms.

**11.17** Write a MATLAB program to implement the quasi-Newton algorithm for general functions. Use the secant method for the line search (e.g., the MATLAB function of Exercise 7.11). Test the various update formulas for  $\mathbf{H}_k$  on Rosenbrock's function (see Exercise 9.4), with an initial condition  $\mathbf{x}^{(0)} = [-2, 2]^\top$ . For this exercise, reinitialize the update direction to the negative gradient every six iterations.

**11.18** Consider the function

$$f(\mathbf{x}) = \frac{x_1^4}{4} + \frac{x_2^2}{2} - x_1 x_2 + x_1 - x_2.$$

- a. Use MATLAB to plot the level sets of  $f$  at levels  $-0.72, -0.6, -0.2, 0.5$ ,  
2. Locate the minimizers of  $f$  from the plots of the level sets.
- b. Apply the DFP algorithm to minimize the function above with the following starting initial conditions: (i)  $[0, 0]^\top$ ; (ii)  $[1.5, 1]^\top$ . Use  $\mathbf{H}_0 = \mathbf{I}_2$ . Does the algorithm converge to the same point for the two initial conditions? If not, explain.

# CHAPTER 12

---

## SOLVING LINEAR EQUATIONS

---

### 12.1 Least-Squares Analysis

Consider a system of linear equations

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \geq n$ , and  $\text{rank } \mathbf{A} = n$ . Note that the number of unknowns,  $n$ , is no larger than the number of equations,  $m$ . If  $\mathbf{b}$  does not belong to the range of  $\mathbf{A}$ , that is, if  $\mathbf{b} \notin \mathcal{R}(\mathbf{A})$ , then this system of equations is said to be *inconsistent* or *overdetermined*. In this case there is no solution to the above set of equations. Our goal then is to find the vector (or vectors)  $\mathbf{x}$  minimizing  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ . This problem is a special case of the nonlinear least-squares problem discussed in Section 9.4.

Let  $\mathbf{x}^*$  be a vector that minimizes  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ ; that is, for all  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 \geq \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2.$$

We refer to the vector  $\mathbf{x}^*$  as a *least-squares solution* to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . In the case where  $\mathbf{A}\mathbf{x} = \mathbf{b}$  has a solution, then the solution is a least-squares solution.

Otherwise, a least-squares solution minimizes the norm of the difference between the left- and right-hand sides of the equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . To characterize least-squares solutions, we need the following lemma.

**Lemma 12.1** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ . Then,  $\text{rank } \mathbf{A} = n$  if and only if  $\text{rank } \mathbf{A}^\top \mathbf{A} = n$  (i.e., the square matrix  $\mathbf{A}^\top \mathbf{A}$  is nonsingular).*  $\square$

*Proof.*  $\Rightarrow$ : Suppose that  $\text{rank } \mathbf{A} = n$ . To show  $\text{rank } \mathbf{A}^\top \mathbf{A} = n$ , it is equivalent to show  $\mathcal{N}(\mathbf{A}^\top \mathbf{A}) = \{\mathbf{0}\}$ . To proceed, let  $\mathbf{x} \in \mathcal{N}(\mathbf{A}^\top \mathbf{A})$ ; that is,  $\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{0}$ . Therefore,

$$\|\mathbf{A}\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{A}^\top \mathbf{A}\mathbf{x} = 0,$$

which implies that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ . Because  $\text{rank } \mathbf{A} = n$ , we have  $\mathbf{x} = \mathbf{0}$ .

$\Leftarrow$ : Suppose that  $\text{rank } \mathbf{A}^\top \mathbf{A} = n$ ; that is,  $\mathcal{N}(\mathbf{A}^\top \mathbf{A}) = \{\mathbf{0}\}$ . To show  $\text{rank } \mathbf{A} = n$ , it is equivalent to show that  $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$ . To proceed, let  $\mathbf{x} \in \mathcal{N}(\mathbf{A})$ ; that is,  $\mathbf{A}\mathbf{x} = \mathbf{0}$ . Then,  $\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{0}$ , and hence  $\mathbf{x} = \mathbf{0}$ .  $\blacksquare$

Recall that we assume throughout that  $\text{rank } \mathbf{A} = n$ . By Lemma 12.1 we conclude that  $(\mathbf{A}^\top \mathbf{A})^{-1}$  exists. The following theorem characterizes the least-squares solution.

**Theorem 12.1** *The unique vector  $\mathbf{x}^*$  that minimizes  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  is given by the solution to the equation  $\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{A}^\top \mathbf{b}$ ; that is,  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ .*  $\square$

*Proof.* Let  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ . First observe that

$$\begin{aligned} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 &= \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*) + (\mathbf{A}\mathbf{x}^* - \mathbf{b})\|^2 \\ &= (\mathbf{A}(\mathbf{x} - \mathbf{x}^*) + (\mathbf{A}\mathbf{x}^* - \mathbf{b}))^\top (\mathbf{A}(\mathbf{x} - \mathbf{x}^*) + (\mathbf{A}\mathbf{x}^* - \mathbf{b})) \\ &= \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 + \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2 + 2[\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{A}\mathbf{x}^* - \mathbf{b}). \end{aligned}$$

We now show that the last term in this equation is zero. Indeed, substituting the expression above for  $\mathbf{x}^*$ ,

$$\begin{aligned} [\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{A}\mathbf{x}^* - \mathbf{b}) &= (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{A}^\top [\mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top - \mathbf{I}_n] \mathbf{b} \\ &= (\mathbf{x} - \mathbf{x}^*)^\top [(\mathbf{A}^\top \mathbf{A})(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top - \mathbf{A}^\top] \mathbf{b} \\ &= (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{A}^\top - \mathbf{A}^\top) \mathbf{b} \\ &= 0. \end{aligned}$$

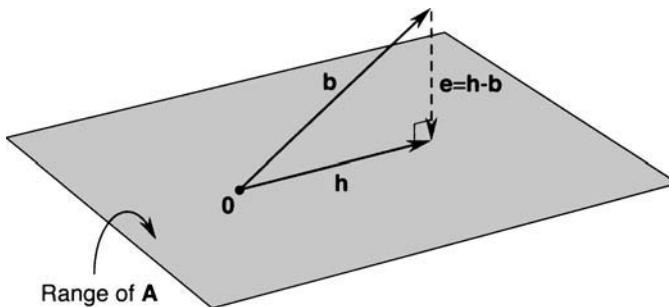
Hence,

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 + \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2.$$

If  $\mathbf{x} \neq \mathbf{x}^*$ , then  $\|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 > 0$ , because  $\text{rank } \mathbf{A} = n$ . Thus, if  $\mathbf{x} \neq \mathbf{x}^*$ , we have

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 > \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2.$$

Thus,  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$  is the unique minimizer of  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ .  $\blacksquare$



**Figure 12.1** Orthogonal projection of  $\mathbf{b}$  on the subspace  $\mathcal{R}(\mathbf{A})$ .

We now give a geometric interpretation of the Theorem 12.1. First note that the columns of  $\mathbf{A}$  span the range  $\mathcal{R}(\mathbf{A})$  of  $\mathbf{A}$ , which is an  $n$ -dimensional subspace of  $\mathbb{R}^m$ . The equation  $\mathbf{Ax} = \mathbf{b}$  has a solution if and only if  $\mathbf{b}$  lies in this  $n$ -dimensional subspace  $\mathcal{R}(\mathbf{A})$ . If  $m = n$ , then  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$  always, and the solution is  $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$ . Suppose now that  $m > n$ . Intuitively, we would expect the “likelihood” of  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$  to be small, because the subspace spanned by the columns of  $\mathbf{A}$  is very “thin.” Therefore, let us suppose that  $\mathbf{b}$  does not belong to  $\mathcal{R}(\mathbf{A})$ . We wish to find a point  $\mathbf{h} \in \mathcal{R}(\mathbf{A})$  that is “closest” to  $\mathbf{b}$ . Geometrically, the point  $\mathbf{h}$  should be such that the vector  $\mathbf{e} = \mathbf{b} - \mathbf{h}$  is orthogonal to the subspace  $\mathcal{R}(\mathbf{A})$  (see Figure 12.1). Recall that a vector  $\mathbf{e} \in \mathbb{R}^m$  is said to be orthogonal to the subspace  $\mathcal{R}(\mathbf{A})$  if it is orthogonal to every vector in this subspace. We call  $\mathbf{h}$  the *orthogonal projection* of  $\mathbf{b}$  onto the subspace  $\mathcal{R}(\mathbf{A})$ . It turns out that  $\mathbf{h} = \mathbf{Ax}^* = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1}\mathbf{A}^\top \mathbf{b}$ . Hence, the vector  $\mathbf{h} \in \mathcal{R}(\mathbf{A})$  minimizing  $\|\mathbf{b} - \mathbf{h}\|$  is exactly the orthogonal projection of  $\mathbf{b}$  onto  $\mathcal{R}(\mathbf{A})$ . In other words, the vector  $\mathbf{x}^*$  minimizing  $\|\mathbf{Ax} - \mathbf{b}\|$  is exactly the vector that makes  $\mathbf{Ax} - \mathbf{b}$  orthogonal to  $\mathcal{R}(\mathbf{A})$ .

To proceed further, we write  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ , where  $\mathbf{a}_1, \dots, \mathbf{a}_n$  are the columns of  $\mathbf{A}$ . The vector  $\mathbf{e}$  is orthogonal to  $\mathcal{R}(\mathbf{A})$  if and only if it is orthogonal to each of the columns  $\mathbf{a}_1, \dots, \mathbf{a}_n$  of  $\mathbf{A}$ . To see this, note that

$$\langle \mathbf{e}, \mathbf{a}_i \rangle = 0, \quad i = 1, \dots, n$$

if and only if for any set of scalars  $\{x_1, x_2, \dots, x_n\}$ , we also have

$$\langle \mathbf{e}, x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n \rangle = 0.$$

Any vector in  $\mathcal{R}(\mathbf{A})$  has the form  $x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n$ .

**Proposition 12.1** Let  $\mathbf{h} \in \mathcal{R}(\mathbf{A})$  be such that  $\mathbf{h} - \mathbf{b}$  is orthogonal to  $\mathcal{R}(\mathbf{A})$ . Then,  $\mathbf{h} = \mathbf{Ax}^* = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1}\mathbf{A}^\top \mathbf{b}$ .  $\square$

*Proof.* Because  $\mathbf{h} \in \mathcal{R}(\mathbf{A}) = \text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n]$ , it has the form  $\mathbf{h} = x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n$ , where  $x_1, \dots, x_n \in \mathbb{R}$ . To find  $x_1, \dots, x_n$ , we use the assumption

that  $\mathbf{e} = \mathbf{h} - \mathbf{b}$  is orthogonal to  $\text{span}[\mathbf{a}_1, \dots, \mathbf{a}_n]$ ; that is, for all  $i = 1, \dots, n$ , we have

$$\langle \mathbf{h} - \mathbf{b}, \mathbf{a}_i \rangle = 0,$$

or, equivalently,

$$\langle \mathbf{h}, \mathbf{a}_i \rangle = \langle \mathbf{b}, \mathbf{a}_i \rangle.$$

Substituting  $\mathbf{h}$  into the equations above, we obtain a set of  $n$  linear equations of the form

$$\langle \mathbf{a}_1, \mathbf{a}_i \rangle x_1 + \cdots + \langle \mathbf{a}_n, \mathbf{a}_i \rangle x_n = \langle \mathbf{b}, \mathbf{a}_i \rangle, \quad i = 1, \dots, n.$$

In matrix notation this system of  $n$  equations can be represented as

$$\begin{bmatrix} \langle \mathbf{a}_1, \mathbf{a}_1 \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_1 \rangle \\ \vdots & & \vdots \\ \langle \mathbf{a}_1, \mathbf{a}_n \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_n \rangle \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \langle \mathbf{b}, \mathbf{a}_1 \rangle \\ \vdots \\ \langle \mathbf{b}, \mathbf{a}_n \rangle \end{bmatrix}.$$

Note that we can write

$$\begin{bmatrix} \langle \mathbf{a}_1, \mathbf{a}_1 \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_1 \rangle \\ \vdots & & \vdots \\ \langle \mathbf{a}_1, \mathbf{a}_n \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_n \rangle \end{bmatrix} = \mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_n^\top \end{bmatrix} [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n].$$

We also note that

$$\begin{bmatrix} \langle \mathbf{b}, \mathbf{a}_1 \rangle \\ \vdots \\ \langle \mathbf{b}, \mathbf{a}_n \rangle \end{bmatrix} = \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_n^\top \end{bmatrix} \mathbf{b}.$$

Because  $\text{rank } \mathbf{A} = n$ ,  $\mathbf{A}^\top \mathbf{A}$  is nonsingular, and thus we conclude that

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \mathbf{x}^*.$$

■

Notice that the matrix

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \langle \mathbf{a}_1, \mathbf{a}_1 \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_1 \rangle \\ \vdots & & \vdots \\ \langle \mathbf{a}_1, \mathbf{a}_n \rangle & \cdots & \langle \mathbf{a}_n, \mathbf{a}_n \rangle \end{bmatrix}$$

plays an important role in the least-squares solution. This matrix is often called the *Gram matrix* (or *Gramian*).

An alternative method of arriving at the least-squares solution is to proceed as follows. First, we write

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{Ax} - \mathbf{b}\|^2 \\ &= (\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b}) \\ &= \frac{1}{2} \mathbf{x}^\top (2\mathbf{A}^\top \mathbf{A}) \mathbf{x} - \mathbf{x}^\top (2\mathbf{A}^\top \mathbf{b}) + \mathbf{b}^\top \mathbf{b}. \end{aligned}$$

Therefore,  $f$  is a quadratic function. The quadratic term is positive definite because  $\text{rank } \mathbf{A} = n$ . Thus, the unique minimizer of  $f$  is obtained by solving the FONC (see Exercise 6.33); that is,

$$\nabla f(\mathbf{x}) = 2\mathbf{A}^\top \mathbf{Ax} - 2\mathbf{A}^\top \mathbf{b} = \mathbf{0}.$$

The only solution to the equation  $\nabla f(\mathbf{x}) = \mathbf{0}$  is  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ .

**Example 12.1** Suppose that you are given two different types of concrete. The first type contains 30% cement, 40% gravel, and 30% sand (all percentages of weight). The second type contains 10% cement, 20% gravel, and 70% sand. How many pounds of each type of concrete should you mix together so that you get a concrete mixture that has as close as possible to a total of 5 pounds of cement, 3 pounds of gravel, and 4 pounds of sand?

The problem can be formulated as a least-squares problem with

$$\mathbf{A} = \begin{bmatrix} 0.3 & 0.1 \\ 0.4 & 0.2 \\ 0.3 & 0.7 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 5 \\ 3 \\ 4 \end{bmatrix},$$

where the decision variable is  $\mathbf{x} = [x_1, x_2]^\top$  and  $x_1$  and  $x_2$  are the amounts of concrete of the first and second types, respectively. After some algebra, we obtain the solution:

$$\begin{aligned} \mathbf{x}^* &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \frac{1}{(0.34)(0.54) - (0.32)^2} \begin{bmatrix} 0.54 & -0.32 \\ -0.32 & 0.34 \end{bmatrix} \begin{bmatrix} 3.9 \\ 3.9 \end{bmatrix} \\ &= \begin{bmatrix} 10.6 \\ 0.961 \end{bmatrix}. \end{aligned}$$

(For a variation of this problem solved using a different method, see Example 15.7.) ■

We now give an example in which least-squares analysis is used to fit measurements by a straight line.

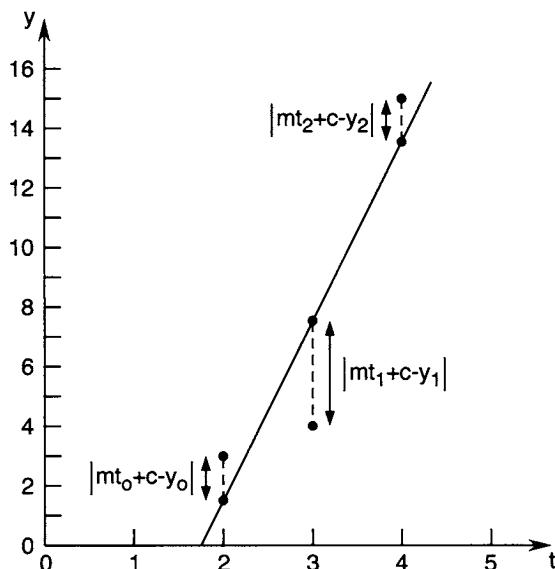
**Table 12.1** Experimental Data for Example 12.2.

$i$	0	1	2
$t_i$	2	3	4
$y_i$	3	4	15

**Example 12.2 Line Fitting.** Suppose that a process has a single input  $t \in \mathbb{R}$  and a single output  $y \in \mathbb{R}$ . Suppose that we perform an experiment on the process, resulting in a number of measurements, as displayed in Table 12.1. The  $i$ th measurement results in the input labeled  $t_i$  and the output labeled  $y_i$ . We would like to find a straight line given by

$$y = mt + c$$

that fits the experimental data. In other words, we wish to find two numbers,  $m$  and  $c$ , such that  $y_i = mt_i + c$ ,  $i = 0, 1, 2$ . However, it is apparent that there is no choice of  $m$  and  $c$  that results in the requirement above; that is, there is no straight line that passes through all three points simultaneously. Therefore, we would like to find the values of  $m$  and  $c$  that best fit the data. A graphical illustration of our problem is shown in Figure 12.2.

**Figure 12.2** Fitting a straight line to experimental data.

We can represent our problem as a system of three linear equations of the form

$$\begin{aligned} 2m + c &= 3 \\ 3m + c &= 4 \\ 4m + c &= 15. \end{aligned}$$

We can write this system of equations as

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 3 & 1 \\ 4 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 4 \\ 15 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} m \\ c \end{bmatrix}.$$

Note that since

$$\text{rank } \mathbf{A} < \text{rank } [\mathbf{A}, \mathbf{b}],$$

the vector  $\mathbf{b}$  does not belong to the range of  $\mathbf{A}$ . Thus, as we have noted before, the system of equations above is inconsistent.

The straight line of best fit is the one that minimizes

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \sum_{i=0}^2 (mt_i + c - y_i)^2.$$

Therefore, our problem lies in the class of least-squares problems. Note that the foregoing function of  $m$  and  $c$  is simply the total squared vertical distance (squared error) between the straight line defined by  $m$  and  $c$  and the experimental points. The solution to our least-squares problem is

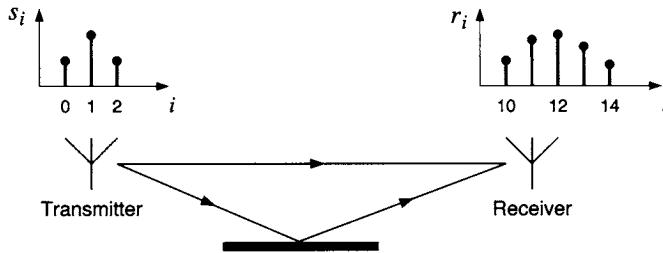
$$\mathbf{x}^* = \begin{bmatrix} m^* \\ c^* \end{bmatrix} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} = \begin{bmatrix} 6 \\ -32/3 \end{bmatrix}.$$

Note that the error vector  $\mathbf{e} = \mathbf{A}\mathbf{x}^* - \mathbf{b}$  is orthogonal to each column of  $\mathbf{A}$ . ■

Next, we give an example of the use of least-squares in wireless communications.

**Example 12.3 Attenuation Estimation.** A wireless transmitter sends a discrete-time signal  $\{s_0, s_1, s_2\}$  (of duration 3) to a receiver, as shown in Figure 12.3. The real number  $s_i$  is the value of the signal at time  $i$ .

The transmitted signal takes two paths to the receiver: a direct path, with delay 10 and attenuation factor  $a_1$ , and an indirect (reflected) path, with delay 12 and attenuation factor  $a_2$ . The received signal is the sum of the signals from these two paths, with their respective delays and attenuation factors.



**Figure 12.3** Wireless transmission in Example 12.3.

Suppose that the received signal is measured from times 10 through 14 as  $r_{10}, r_{11}, \dots, r_{14}$ , as shown in the figure. We wish to compute the least-squares estimates of  $a_1$  and  $a_2$ , based on the following values:

$$\frac{s_0 \quad s_1 \quad s_2 \quad r_{10} \quad r_{11} \quad r_{12} \quad r_{13} \quad r_{14}}{1 \quad 2 \quad 1 \quad 4 \quad 7 \quad 8 \quad 6 \quad 3}.$$

The problem can be posed as a least-squares problem with

$$\mathbf{A} = \begin{bmatrix} s_0 & 0 \\ s_1 & 0 \\ s_2 & s_0 \\ 0 & s_1 \\ 0 & s_2 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} r_{10} \\ r_{11} \\ r_{12} \\ r_{13} \\ r_{14} \end{bmatrix}.$$

The least-squares estimate is given by

$$\begin{aligned} \begin{bmatrix} a_1^* \\ a_2^* \end{bmatrix} &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \begin{bmatrix} \|\mathbf{s}\|^2 & s_0 s_2 \\ s_0 s_2 & \|\mathbf{s}\|^2 \end{bmatrix}^{-1} \begin{bmatrix} s_0 r_{10} + s_1 r_{11} + s_2 r_{12} \\ s_0 r_{12} + s_1 r_{13} + s_2 r_{14} \end{bmatrix} \\ &= \begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix}^{-1} \begin{bmatrix} 4 + 14 + 8 \\ 8 + 12 + 3 \end{bmatrix} \\ &= \frac{1}{35} \begin{bmatrix} 6 & -1 \\ -1 & 6 \end{bmatrix} \begin{bmatrix} 26 \\ 23 \end{bmatrix} \\ &= \frac{1}{35} \begin{bmatrix} 133 \\ 112 \end{bmatrix}. \end{aligned}$$

■

We now give a simple example where the least-squares method is used in digital signal processing.

**Example 12.4 Discrete Fourier Series.** Suppose that we are given a discrete-time signal, represented by the vector

$$\mathbf{b} = [b_1, b_2, \dots, b_m]^\top.$$

We wish to approximate this signal by a sum of sinusoids. Specifically, we approximate  $\mathbf{b}$  by the vector

$$y_0 \mathbf{c}^{(0)} + \sum_{k=1}^n \left( y_k \mathbf{c}^{(k)} + z_k \mathbf{s}^{(k)} \right),$$

where  $y_0, y_1, \dots, y_n, z_1, \dots, z_n \in \mathbb{R}$  and the vectors  $\mathbf{c}^{(k)}$  and  $\mathbf{s}^{(k)}$  are given by

$$\begin{aligned} \mathbf{c}^{(0)} &= \left[ \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, \dots, \frac{1}{\sqrt{2}} \right]^\top, \\ \mathbf{c}^{(k)} &= \left[ \cos \left( 1 \frac{2k\pi}{m} \right), \cos \left( 2 \frac{2k\pi}{m} \right), \dots, \cos \left( m \frac{2k\pi}{m} \right) \right]^\top, \quad k = 1, \dots, n, \\ \mathbf{s}^{(k)} &= \left[ \sin \left( 1 \frac{2k\pi}{m} \right), \sin \left( 2 \frac{2k\pi}{m} \right), \dots, \sin \left( m \frac{2k\pi}{m} \right) \right]^\top, \quad k = 1, \dots, n. \end{aligned}$$

We call the sum of sinusoids above a *discrete Fourier series* (although, strictly speaking, it is not a series but a finite sum). We wish to find  $y_0, y_1, \dots, y_n, z_1, \dots, z_n$  such that

$$\left\| \left( y_0 \mathbf{c}^{(0)} + \sum_{k=1}^n y_k \mathbf{c}^{(k)} + z_k \mathbf{s}^{(k)} \right) - \mathbf{b} \right\|^2$$

is minimized.

To proceed, we define

$$\begin{aligned} \mathbf{A} &= \left[ \mathbf{c}^{(0)}, \mathbf{c}^{(1)}, \dots, \mathbf{c}^{(n)}, \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(n)} \right], \\ \mathbf{x} &= [y_0, y_1, \dots, y_n, z_1, \dots, z_n]^\top. \end{aligned}$$

Our problem can be reformulated as minimizing

$$\|\mathbf{Ax} - \mathbf{b}\|^2.$$

We assume that  $m \geq 2n+1$ . To find the solution, we first compute  $\mathbf{A}^\top \mathbf{A}$ . We make use of the following trigonometric identities: For any nonzero integer  $k$  that is not an integral multiple of  $m$ , we have

$$\sum_{i=1}^m \cos \left( i \frac{2k\pi}{m} \right) = 0,$$

$$\sum_{i=1}^m \sin \left( i \frac{2k\pi}{m} \right) = 0.$$

With the aid of these identities, we can verify that

$$\begin{aligned}\mathbf{c}^{(k)\top} \mathbf{c}^{(j)} &= \begin{cases} m/2 & \text{if } k = j \\ 0 & \text{otherwise} \end{cases} \\ \mathbf{s}^{(k)\top} \mathbf{s}^{(j)} &= \begin{cases} m/2 & \text{if } k = j \\ 0 & \text{otherwise} \end{cases} \\ \mathbf{c}^{(k)\top} \mathbf{s}^{(j)} &= 0 \quad \text{for any } k, j.\end{aligned}$$

Hence,

$$\mathbf{A}^\top \mathbf{A} = \frac{m}{2} \mathbf{I}_{2n+1},$$

which is clearly nonsingular, with inverse

$$(\mathbf{A}^\top \mathbf{A})^{-1} = \frac{2}{m} \mathbf{I}_{2n+1}.$$

Therefore, the solution to our problem is

$$\begin{aligned}\mathbf{x}^* &= [y_0^*, y_1^*, \dots, y_n^*, z_1^*, \dots, z_n^*]^\top \\ &= (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b} \\ &= \frac{2}{m} \mathbf{A}^\top \mathbf{b}.\end{aligned}$$

We represent the solution as

$$\begin{aligned}y_0^* &= \frac{\sqrt{2}}{m} \sum_{i=1}^m b_i, \\ y_k^* &= \frac{2}{m} \sum_{i=1}^m b_i \cos\left(i \frac{2k\pi}{m}\right), \quad k = 1, \dots, n, \\ z_k^* &= \frac{2}{m} \sum_{i=1}^m b_i \sin\left(i \frac{2k\pi}{m}\right), \quad k = 1, \dots, n.\end{aligned}$$

We call these *discrete Fourier coefficients*. ■

Finally, we show how least-squares analysis can be used to derive formulas for orthogonal projectors.

**Example 12.5 Orthogonal Projectors.** Let  $\mathcal{V} \subset \mathbb{R}^n$  be a subspace. Given a vector  $\mathbf{x} \in \mathbb{R}^n$ , we write the orthogonal decomposition of  $\mathbf{x}$  as

$$\mathbf{x} = \mathbf{x}_\mathcal{V} + \mathbf{x}_{\mathcal{V}^\perp},$$

where  $\mathbf{x}_\mathcal{V} \in \mathcal{V}$  is the orthogonal projection of  $\mathbf{x}$  onto  $\mathcal{V}$  and  $\mathbf{x}_{\mathcal{V}^\perp} \in \mathcal{V}^\perp$  is the orthogonal projection of  $\mathbf{x}$  onto  $\mathcal{V}^\perp$ . (See Section 3.3; also recall that  $\mathcal{V}^\perp$  is

the orthogonal complement of  $\mathcal{V}$ .) We can write  $\mathbf{x}_\mathcal{V} = \mathbf{P}\mathbf{x}$  for some matrix  $\mathbf{P}$  called the *orthogonal projector*. In the following, we derive expressions for  $\mathbf{P}$  for the case where  $\mathcal{V} = \mathcal{R}(\mathbf{A})$  and the case where  $\mathcal{V} = \mathcal{N}(\mathbf{A})$ .

Consider a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ , and  $\text{rank } \mathbf{A} = n$ . Let  $\mathcal{V} = \mathcal{R}(\mathbf{A})$  be the range of  $\mathbf{A}$  (note that any subspace can be written as the range of some matrix). In this case we can write an expression for  $\mathbf{P}$  in terms of  $\mathbf{A}$ , as follows. By Proposition 12.1 we have  $\mathbf{x}_\mathcal{V} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{x}$ , whence  $\mathbf{P} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$ . Note that by Proposition 12.1, we may also write

$$\mathbf{x}_\mathcal{V} = \arg \min_{\mathbf{y} \in \mathcal{V}} \|\mathbf{y} - \mathbf{x}\|.$$

Next, consider a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . Let  $\mathcal{V} = \mathcal{N}(\mathbf{A})$  be the nullspace of  $\mathbf{A}$  (note that any subspace can be written as the nullspace of some matrix). To derive an expression for the orthogonal projector  $\mathbf{P}$  in terms of  $\mathbf{A}$  for this case, we use the formula derived above and the identity  $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top)$  (see Theorem 3.4). Indeed, if  $\mathcal{U} = \mathcal{R}(\mathbf{A}^\top)$ , then the orthogonal decomposition with respect to  $\mathcal{U}$  is  $\mathbf{x} = \mathbf{x}_\mathcal{U} + \mathbf{x}_{\mathcal{U}^\perp}$ , where  $\mathbf{x}_\mathcal{U} = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}$  (using the formula derived above). Because  $\mathcal{N}(\mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top)$ , we deduce that  $\mathbf{x}_{\mathcal{V}^\perp} = \mathbf{x}_\mathcal{U} = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x}$ . Hence,

$$\mathbf{x}_\mathcal{V} = \mathbf{x} - \mathbf{x}_{\mathcal{V}^\perp} = \mathbf{x} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{x} = (\mathbf{I} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A})\mathbf{x}.$$

Thus, the orthogonal projector in this case is  $\mathbf{P} = \mathbf{I} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}$ . ■

## 12.2 The Recursive Least-Squares Algorithm

Consider again the example in Section 12.1. We are given experimental points  $(t_0, y_0)$ ,  $(t_1, y_1)$ , and  $(t_2, y_2)$ , and we find the parameters  $m^*$  and  $c^*$  of the straight line that best fits these data in the least-squares sense. Suppose that we are now given an extra measurement point  $(t_3, y_3)$ , so that we now have a set of four experimental data points:  $(t_0, y_0)$ ,  $(t_1, y_1)$ ,  $(t_2, y_2)$ , and  $(t_3, y_3)$ . We can similarly go through the procedure for finding the parameters of the line of best fit for this set of four points. However, as we shall see, there is a more efficient way: We can use previous calculations of  $m^*$  and  $c^*$  for the three data points to calculate the parameters for the four data points. In effect, we simply “update” our values of  $m^*$  and  $c^*$  to accommodate the new data point. This procedure, called the *recursive least-squares (RLS) algorithm*, is the topic of this section.

To derive the RLS algorithm, first consider the problem of minimizing  $\|\mathbf{A}_0 \mathbf{x} - \mathbf{b}^{(0)}\|^2$ . We know that the solution to this is given by  $\mathbf{x}^{(0)} = \mathbf{G}_0^{-1} \mathbf{A}_0^\top \mathbf{b}^{(0)}$ , where  $\mathbf{G}_0 = \mathbf{A}_0^\top \mathbf{A}_0$ . Suppose now that we are given more data, in the form of a matrix  $\mathbf{A}_1$  and a vector  $\mathbf{b}^{(1)}$ . Consider now the problem of minimizing

$$\left\| \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \end{bmatrix} \right\|^2.$$

The solution is given by

$$\mathbf{x}^{(1)} = \mathbf{G}_1^{-1} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \end{bmatrix},$$

where

$$\mathbf{G}_1 = \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}.$$

Our goal is to write  $\mathbf{x}^{(1)}$  as a function of  $\mathbf{x}^{(0)}$ ,  $\mathbf{G}_0$ , and the new data  $\mathbf{A}_1$  and  $\mathbf{b}^{(1)}$ . To this end, we first write  $\mathbf{G}_1$  as

$$\begin{aligned} \mathbf{G}_1 &= [\mathbf{A}_0^\top \ \mathbf{A}_1^\top] \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix} \\ &= \mathbf{A}_0^\top \mathbf{A}_0 + \mathbf{A}_1^\top \mathbf{A}_1 \\ &= \mathbf{G}_0 + \mathbf{A}_1^\top \mathbf{A}_1. \end{aligned}$$

Next, we write

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \end{bmatrix} &= [\mathbf{A}_0^\top \ \mathbf{A}_1^\top] \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \end{bmatrix} \\ &= \mathbf{A}_0^\top \mathbf{b}^{(0)} + \mathbf{A}_1^\top \mathbf{b}^{(1)}. \end{aligned}$$

To proceed further, we write  $\mathbf{A}_0^\top \mathbf{b}^{(0)}$  as

$$\begin{aligned} \mathbf{A}_0^\top \mathbf{b}^{(0)} &= \mathbf{G}_0 \mathbf{G}_0^{-1} \mathbf{A}_0^\top \mathbf{b}^{(0)} \\ &= \mathbf{G}_0 \mathbf{x}^{(0)} \\ &= (\mathbf{G}_1 - \mathbf{A}_1^\top \mathbf{A}_1) \mathbf{x}^{(0)} \\ &= \mathbf{G}_1 \mathbf{x}^{(0)} - \mathbf{A}_1^\top \mathbf{A}_1 \mathbf{x}^{(0)}. \end{aligned}$$

Combining these formulas, we see that we can write  $\mathbf{x}^{(1)}$  as

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{G}_1^{-1} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \end{bmatrix} \\ &= \mathbf{G}_1^{-1} (\mathbf{G}_1 \mathbf{x}^{(0)} - \mathbf{A}_1^\top \mathbf{A}_1 \mathbf{x}^{(0)} + \mathbf{A}_1^\top \mathbf{b}^{(1)}) \\ &= \mathbf{x}^{(0)} + \mathbf{G}_1^{-1} \mathbf{A}_1^\top (\mathbf{b}^{(1)} - \mathbf{A}_1 \mathbf{x}^{(0)}), \end{aligned}$$

where  $\mathbf{G}_1$  can be calculated using

$$\mathbf{G}_1 = \mathbf{G}_0 + \mathbf{A}_1^\top \mathbf{A}_1.$$

We note that with this formula,  $\mathbf{x}^{(1)}$  can be computed using only  $\mathbf{x}^{(0)}$ ,  $\mathbf{A}_1$ ,  $\mathbf{b}^{(1)}$ , and  $\mathbf{G}_0$ . Hence, we have a way of using our previous efforts in calculating  $\mathbf{x}^{(0)}$  to compute  $\mathbf{x}^{(1)}$ , without having to compute  $\mathbf{x}^{(1)}$  directly from scratch. The solution  $\mathbf{x}^{(1)}$  is obtained from  $\mathbf{x}^{(0)}$  by a simple update equation that adds to  $\mathbf{x}^{(0)}$  a “correction term”  $\mathbf{G}_1^{-1}\mathbf{A}_1^\top (\mathbf{b}^{(1)} - \mathbf{A}_1\mathbf{x}^{(0)})$ . Observe that if the new data are consistent with the old data, that is,  $\mathbf{A}_1\mathbf{x}^{(0)} = \mathbf{b}^{(1)}$ , then the correction term is 0 and the updated solution  $\mathbf{x}^{(1)}$  is equal to the previous solution  $\mathbf{x}^{(0)}$ .

We can generalize the argument above to write a recursive algorithm for updating the least-squares solution as new data arrive. At the  $(k+1)$ th iteration, we have

$$\begin{aligned}\mathbf{G}_{k+1} &= \mathbf{G}_k + \mathbf{A}_{k+1}^\top \mathbf{A}_{k+1} \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \mathbf{G}_{k+1}^{-1} \mathbf{A}_{k+1}^\top (\mathbf{b}^{(k+1)} - \mathbf{A}_{k+1}\mathbf{x}^{(k)}).\end{aligned}$$

The vector  $\mathbf{b}^{(k+1)} - \mathbf{A}_{k+1}\mathbf{x}^{(k)}$  is often called the *innovation*. As before, observe that if the innovation is zero, then the updated solution  $\mathbf{x}^{(k+1)}$  is equal to the previous solution  $\mathbf{x}^{(k)}$ .

We can see from the above that to compute  $\mathbf{x}^{(k+1)}$  from  $\mathbf{x}^{(k)}$  we need  $\mathbf{G}_{k+1}^{-1}$ , rather than  $\mathbf{G}_{k+1}$ . It turns out that we can derive an update formula for  $\mathbf{G}_{k+1}^{-1}$  itself. To do so, we need the following technical lemma, which is a generalization of the Sherman-Morrison formula (Lemma 11.1), due to Woodbury, and hence also called the *Sherman-Morrison-Woodbury formula* (see [63, p. 124] or [53, p. 50]).

**Lemma 12.2** *Let  $\mathbf{A}$  be a nonsingular matrix. Let  $\mathbf{U}$  and  $\mathbf{V}$  be matrices such that  $\mathbf{I} + \mathbf{V}\mathbf{A}^{-1}\mathbf{U}$  is nonsingular. Then,  $\mathbf{A} + \mathbf{U}\mathbf{V}$  is nonsingular, and*

$$(\mathbf{A} + \mathbf{U}\mathbf{V})^{-1} = \mathbf{A}^{-1} - (\mathbf{A}^{-1}\mathbf{U})(\mathbf{I} + \mathbf{V}\mathbf{A}^{-1}\mathbf{U})^{-1}(\mathbf{V}\mathbf{A}^{-1}).$$

□

*Proof.* We can prove the result easily by verification. ■

Using Lemma 12.2 we get

$$\begin{aligned}\mathbf{G}_{k+1}^{-1} &= \left( \mathbf{G}_k + \mathbf{A}_{k+1}^\top \mathbf{A}_{k+1} \right)^{-1} \\ &= \mathbf{G}_k^{-1} - \mathbf{G}_k^{-1} \mathbf{A}_{k+1}^\top (\mathbf{I} + \mathbf{A}_{k+1}\mathbf{G}_k^{-1}\mathbf{A}_{k+1}^\top)^{-1} \mathbf{A}_{k+1}\mathbf{G}_k^{-1}.\end{aligned}$$

For simplicity of notation, we rewrite  $\mathbf{G}_k^{-1}$  as  $\mathbf{P}_k$ .

We summarize by writing the RLS algorithm using  $\mathbf{P}_k$ :

$$\begin{aligned}\mathbf{P}_{k+1} &= \mathbf{P}_k - \mathbf{P}_k \mathbf{A}_{k+1}^\top (\mathbf{I} + \mathbf{A}_{k+1}\mathbf{P}_k\mathbf{A}_{k+1}^\top)^{-1} \mathbf{A}_{k+1}\mathbf{P}_k, \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \mathbf{P}_{k+1} \mathbf{A}_{k+1}^\top (\mathbf{b}^{(k+1)} - \mathbf{A}_{k+1}\mathbf{x}^{(k)}).\end{aligned}$$

In the special case where the new data at each step are such that  $\mathbf{A}_{k+1}$  is a matrix consisting of a single row,  $\mathbf{A}_{k+1} = \mathbf{a}_{k+1}^\top$ , and  $\mathbf{b}^{(k+1)}$  is a scalar,  $\mathbf{b}^{(k+1)} = b_{k+1}$ , we get

$$\begin{aligned}\mathbf{P}_{k+1} &= \mathbf{P}_k - \frac{\mathbf{P}_k \mathbf{a}_{k+1} \mathbf{a}_{k+1}^\top \mathbf{P}_k}{1 + \mathbf{a}_{k+1}^\top \mathbf{P}_k \mathbf{a}_{k+1}}, \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \mathbf{P}_{k+1} \mathbf{a}_{k+1} \left( b_{k+1} - \mathbf{a}_{k+1}^\top \mathbf{x}^{(k)} \right).\end{aligned}$$

**Example 12.6** Let

$$\begin{aligned}\mathbf{A}_0 &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, & \mathbf{b}^{(0)} &= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \\ \mathbf{A}_1 &= \mathbf{a}_1^\top = [2 \ 1], & \mathbf{b}^{(1)} &= b_1 = [3], \\ \mathbf{A}_2 &= \mathbf{a}_2^\top = [3 \ 1], & \mathbf{b}^{(2)} &= b_2 = [4].\end{aligned}$$

First compute the vector  $\mathbf{x}^{(0)}$  minimizing  $\|\mathbf{A}_0 \mathbf{x} - \mathbf{b}^{(0)}\|^2$ . Then, use the RLS algorithm to find  $\mathbf{x}^{(2)}$  minimizing

$$\left\| \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \end{bmatrix} \right\|^2.$$

We have

$$\begin{aligned}\mathbf{P}_0 &= (\mathbf{A}_0^\top \mathbf{A}_0)^{-1} = \begin{bmatrix} 2/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix}, \\ \mathbf{x}^{(0)} &= \mathbf{P}_0 \mathbf{A}_0^\top \mathbf{b}^{(0)} = \begin{bmatrix} 2/3 \\ 2/3 \end{bmatrix}.\end{aligned}$$

Applying the RLS algorithm twice, we get

$$\begin{aligned}\mathbf{P}_1 &= \mathbf{P}_0 - \frac{\mathbf{P}_0 \mathbf{a}_1 \mathbf{a}_1^\top \mathbf{P}_0}{1 + \mathbf{a}_1^\top \mathbf{P}_0 \mathbf{a}_1} = \begin{bmatrix} 1/3 & -1/3 \\ -1/3 & 2/3 \end{bmatrix}, \\ \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} + \mathbf{P}_1 \mathbf{a}_1 \left( b_1 - \mathbf{a}_1^\top \mathbf{x}^{(0)} \right) = \begin{bmatrix} 1 \\ 2/3 \end{bmatrix}, \\ \mathbf{P}_2 &= \mathbf{P}_1 - \frac{\mathbf{P}_1 \mathbf{a}_2 \mathbf{a}_2^\top \mathbf{P}_1}{1 + \mathbf{a}_2^\top \mathbf{P}_1 \mathbf{a}_2} = \begin{bmatrix} 1/6 & -1/4 \\ -1/4 & 5/8 \end{bmatrix}, \\ \mathbf{x}^{(2)} &= \mathbf{x}^{(1)} + \mathbf{P}_2 \mathbf{a}_2 \left( b_2 - \mathbf{a}_2^\top \mathbf{x}^{(1)} \right) = \begin{bmatrix} 13/12 \\ 5/8 \end{bmatrix}.\end{aligned}$$

We can easily check our solution by computing  $\mathbf{x}^{(2)}$  directly using the formula  $\mathbf{x}^{(2)} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$ , where

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{A}_1 \\ \mathbf{A}_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}^{(0)} \\ \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \end{bmatrix}.$$

■

### 12.3 Solution to a Linear Equation with Minimum Norm

Consider now a system of linear equations

$$\mathbf{Ax} = \mathbf{b},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . Note that the number of equations is no larger than the number of unknowns. There may exist an infinite number of solutions to this system of equations. However, as we shall see, there is only one solution that is closest to the origin: the solution to  $\mathbf{Ax} = \mathbf{b}$  whose norm  $\|\mathbf{x}\|$  is minimal. Let  $\mathbf{x}^*$  be this solution; that is,  $\mathbf{Ax}^* = \mathbf{b}$  and  $\|\mathbf{x}^*\| \leq \|\mathbf{x}\|$  for any  $\mathbf{x}$  such that  $\mathbf{Ax} = \mathbf{b}$ . In other words,  $\mathbf{x}^*$  is the solution to the problem

$$\begin{aligned} &\text{minimize} && \|\mathbf{x}\| \\ &\text{subject to} && \mathbf{Ax} = \mathbf{b}. \end{aligned}$$

In Part IV, we study problems of this type in more detail.

**Theorem 12.2** *The unique solution  $\mathbf{x}^*$  to  $\mathbf{Ax} = \mathbf{b}$  that minimizes the norm  $\|\mathbf{x}\|$  is given by*

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}.$$

□

*Proof.* Let  $\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}$ . Note that

$$\begin{aligned} \|\mathbf{x}\|^2 &= \|(\mathbf{x} - \mathbf{x}^*) + \mathbf{x}^*\|^2 \\ &= ((\mathbf{x} - \mathbf{x}^*) + \mathbf{x}^*)^\top ((\mathbf{x} - \mathbf{x}^*) + \mathbf{x}^*) \\ &= \|\mathbf{x} - \mathbf{x}^*\|^2 + \|\mathbf{x}^*\|^2 + 2\mathbf{x}^{*\top}(\mathbf{x} - \mathbf{x}^*). \end{aligned}$$

We now show that

$$\mathbf{x}^{*\top}(\mathbf{x} - \mathbf{x}^*) = 0.$$

Indeed,

$$\begin{aligned} \mathbf{x}^{*\top}(\mathbf{x} - \mathbf{x}^*) &= [\mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}]^\top [\mathbf{x} - \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}] \\ &= \mathbf{b}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} [\mathbf{A} \mathbf{x} - (\mathbf{A} \mathbf{A}^\top)(\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}] \\ &= \mathbf{b}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} [\mathbf{b} - \mathbf{b}] = 0. \end{aligned}$$

Therefore,

$$\|\mathbf{x}\|^2 = \|\mathbf{x}^*\|^2 + \|\mathbf{x} - \mathbf{x}^*\|^2.$$

Because  $\|\mathbf{x} - \mathbf{x}^*\|^2 > 0$  for all  $\mathbf{x} \neq \mathbf{x}^*$ , it follows that for all  $\mathbf{x} \neq \mathbf{x}^*$ ,

$$\|\mathbf{x}\|^2 > \|\mathbf{x}^*\|^2,$$

which implies that

$$\|\mathbf{x}\| > \|\mathbf{x}^*\|. \quad \blacksquare$$

**Example 12.7** Find the point closest to the origin of  $\mathbb{R}^3$  on the line of intersection of the two planes defined by the following two equations:

$$\begin{aligned} x_1 + 2x_2 - x_3 &= 1, \\ 4x_1 + x_2 + 3x_3 &= 0. \end{aligned}$$

Note that this problem is equivalent to the problem

$$\begin{aligned} &\text{minimize} \quad \|\mathbf{x}\| \\ &\text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -1 \\ 4 & 1 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Thus, the solution to the problem is

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} = \begin{bmatrix} 0.0952 \\ 0.3333 \\ -0.2381 \end{bmatrix}. \quad \blacksquare$$

In the next section we discuss an iterative algorithm for solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

## 12.4 Kaczmarz's Algorithm

As in Section 12.3, let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . We now discuss an iterative algorithm for solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , originally analyzed by Kaczmarz in 1937 [70]. The algorithm converges to the vector  $\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b}$  without explicitly having to invert the matrix  $\mathbf{A}\mathbf{A}^\top$ . This is important from a practical point of view, especially when  $\mathbf{A}$  has many rows.

Let  $\mathbf{a}_j^\top$  denote the  $j$ th row of  $\mathbf{A}$ , and  $b_j$  the  $j$ th component of  $\mathbf{b}$ , and  $\mu$  a positive scalar,  $0 < \mu < 2$ . With this notation, Kaczmarz's algorithm is:

1. Set  $i := 0$ , initial condition  $\mathbf{x}^{(0)}$ .
2. For  $j = 1, \dots, m$ , set  

$$\mathbf{x}^{(im+j)} = \mathbf{x}^{(im+j-1)} + \mu (b_j - \mathbf{a}_j^\top \mathbf{x}^{(im+j-1)}) \frac{\mathbf{a}_j}{\mathbf{a}_j^\top \mathbf{a}_j}.$$
3. Set  $i := i + 1$ ; go to step 2.

In words, Kaczmarz's algorithm works as follows. For the first  $m$  iterations ( $k = 0, \dots, m-1$ ), we have

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu (b_{k+1} - \mathbf{a}_{k+1}^\top \mathbf{x}^{(k)}) \frac{\mathbf{a}_{k+1}}{\mathbf{a}_{k+1}^\top \mathbf{a}_{k+1}},$$

where, in each iteration, we use rows of  $\mathbf{A}$  and corresponding components of  $\mathbf{b}$  successively. For the  $(m+1)$ th iteration, we revert back to the first row of  $\mathbf{A}$  and the first component of  $\mathbf{b}$ ; that is,

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + \mu (b_1 - \mathbf{a}_1^\top \mathbf{x}^{(m)}) \frac{\mathbf{a}_1}{\mathbf{a}_1^\top \mathbf{a}_1}.$$

We continue with the  $(m+2)$ th iteration using the second row of  $\mathbf{A}$  and second component of  $\mathbf{b}$ , and so on, repeating the cycle every  $m$  iteration. We can view the scalar  $\mu$  as the step size of the algorithm. The reason for requiring that  $\mu$  be between 0 and 2 will become apparent from the convergence analysis.

We now prove the convergence of Kaczmarz's algorithm, using ideas from Kaczmarz's original paper [70] and subsequent exposition by Parks [102].

**Theorem 12.3** *In Kaczmarz's algorithm, if  $\mathbf{x}^{(0)} = \mathbf{0}$ , then  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b}$  as  $k \rightarrow \infty$ .*  $\square$

*Proof.* We may assume without loss of generality that  $\|\mathbf{a}_i\| = 1$ ,  $i = 1, \dots, m$ . For if not, we simply replace each  $\mathbf{a}_i$  by  $\mathbf{a}_i/\|\mathbf{a}\|$  and each  $b_i$  by  $b_i/\|\mathbf{a}\|$ .

We first introduce the following notation. For each  $j = 0, 1, 2, \dots$ , let  $R(j)$  denote the unique integer in  $\{0, \dots, m-1\}$  satisfying  $j = lm + R(j)$  for some integer  $l$ ; that is,  $R(j)$  is the remainder that results if we divide  $j$  by  $m$ .

Using the notation above, we can write Kaczmarz's algorithm as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mu (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \mathbf{a}_{R(k)+1}.$$

Using the identity  $\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle$ , we obtain

$$\begin{aligned} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 &= \|\mathbf{x}^{(k)} - \mathbf{x}^* + \mu (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \mathbf{a}_{R(k)+1}\|^2 \\ &= \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 + \mu^2 (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)})^2 \\ &\quad + 2\mu (b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)}) \mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*). \end{aligned}$$

Substituting  $\mathbf{a}_{R(k)+1}^\top \mathbf{x}^* = b_{R(k)+1}$  into this equation, we get

$$\begin{aligned} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 &= \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \mu(2 - \mu)(b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)})^2 \\ &= \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \mu(2 - \mu)(\mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*))^2. \end{aligned}$$

Because  $0 < \mu < 2$ , the second term on the right-hand side is nonnegative, and hence

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2 \leq \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2.$$

Therefore,  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$  is a nonincreasing sequence that is bounded below, because  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 \geq 0$  for all  $k$ . Hence,  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$  converges (see Theorem 5.3). Furthermore, we may write

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 - \mu(2 - \mu) \sum_{i=0}^{k-1} (\mathbf{a}_{R(i)+1}^\top (\mathbf{x}^{(i)} - \mathbf{x}^*))^2.$$

Because  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$  converges, we conclude that

$$\sum_{i=0}^{\infty} (\mathbf{a}_{R(i)+1}^\top (\mathbf{x}^{(i)} - \mathbf{x}^*))^2 < \infty,$$

which implies that

$$\mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*) \rightarrow 0.$$

Observe that

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 = \mu^2 (\mathbf{b}_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k)})^2 = \mu^2 (\mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k)} - \mathbf{x}^*))^2$$

and therefore  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 \rightarrow 0$ . Note also that because  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$  converges,  $\{\mathbf{x}^{(k)}\}$  is a bounded sequence (see Theorem 5.2).

Following Kaczmarz [70], we introduce the notation  $\mathbf{x}^{(r,s)} \triangleq \mathbf{x}^{(rm+s)}$ ,  $r = 0, 1, 2, \dots$ ,  $s = 0, \dots, m-1$ . With this notation, we have, for each  $s = 0, \dots, m-1$ ,

$$\mathbf{a}_{s+1}^\top (\mathbf{x}^{(r,s)} - \mathbf{x}^*) \rightarrow 0$$

as  $r \rightarrow \infty$ . Consider now the sequence  $\{\mathbf{x}^{(r,0)} : r \geq 0\}$ . Because this sequence is bounded, we conclude that it has a convergent subsequence—this follows from the Bolzano-Weierstrass theorem (see [2, p. 70]; see also Section 5.1 for a discussion of sequences and subsequences). Denote this convergent subsequence by  $\{\mathbf{x}^{(r,0)} : r \in \mathcal{E}\}$ , where  $\mathcal{E}$  is a subset of  $\{0, 1, \dots\}$ . Let  $\mathbf{z}^*$  be the limit of  $\{\mathbf{x}^{(r,0)} : r \in \mathcal{E}\}$ . Hence,

$$\mathbf{a}_1^\top (\mathbf{z}^* - \mathbf{x}^*) = 0.$$

Next, note that because  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 \rightarrow 0$  as  $k \rightarrow \infty$ , we also have  $\|\mathbf{x}^{(r,1)} - \mathbf{x}^{(r,0)}\|^2 \rightarrow 0$  as  $r \rightarrow \infty$ . Therefore, the subsequence  $\{\mathbf{x}^{(r,1)} : r \in \mathcal{E}\}$  also converges to  $\mathbf{z}^*$ . Hence,

$$\mathbf{a}_2^\top (\mathbf{z}^* - \mathbf{x}^*) = 0.$$

Repeating the argument, we conclude that for each  $i = 1, \dots, m$ ,

$$\mathbf{a}_i^\top (\mathbf{z}^* - \mathbf{x}^*) = 0.$$

In matrix notation, the equations above take the form

$$\mathbf{A}(\mathbf{z}^* - \mathbf{x}^*) = \mathbf{0}.$$

Now,  $\mathbf{x}^{(k)} \in \mathcal{R}(\mathbf{A}^\top)$  for all  $k$  because  $\mathbf{x}^{(0)} = \mathbf{0}$  (see Exercise 12.25). Therefore,  $\mathbf{z}^* \in \mathcal{R}(\mathbf{A}^\top)$ , because  $\mathcal{R}(\mathbf{A}^\top)$  is closed. Hence, there exists  $\mathbf{y}^*$  such that  $\mathbf{z}^* = \mathbf{A}^\top \mathbf{y}^*$ . Thus,

$$\begin{aligned}\mathbf{A}(\mathbf{z}^* - \mathbf{x}^*) &= \mathbf{A}(\mathbf{A}^\top \mathbf{y}^* - \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}) \\ &= (\mathbf{A} \mathbf{A}^\top) \mathbf{y}^* - \mathbf{b} \\ &= \mathbf{0}.\end{aligned}$$

Because  $\text{rank } \mathbf{A} = m$ ,  $\mathbf{y}^* = (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{b}$  and hence  $\mathbf{z}^* = \mathbf{x}^*$ . Therefore, the subsequence  $\{\|\mathbf{x}^{r,0} - \mathbf{x}^*\|^2 : r \in \mathcal{E}\}$  converges to 0. Because  $\{\|\mathbf{x}^{r,0} - \mathbf{x}^*\|^2 : r \in \mathcal{E}\}$  is a subsequence of the convergent sequence  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$ , we conclude that the sequence  $\{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2\}$  converges to 0; that is,  $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ . ■

For the case where  $\mathbf{x}^{(0)} \neq \mathbf{0}$ , Kaczmarz's algorithm converges to the unique point on  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$  minimizing the distance  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  (see Exercise 12.26).

If we set  $\mu = 1$ , Kaczmarz's algorithm has the property that at each iteration  $k$ , the “error”  $b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k+1)}$  satisfies

$$b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k+1)} = 0$$

(see Exercise 12.28). Substituting  $b_{R(k)+1} = \mathbf{a}_{R(k)+1}^\top \mathbf{x}^*$ , we may write

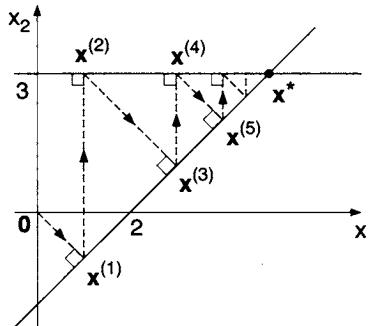
$$\mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k+1)} - \mathbf{x}^*) = 0.$$

Hence, the difference between  $\mathbf{x}^{(k+1)}$  and the solution  $\mathbf{x}^*$  is orthogonal to  $\mathbf{a}_{R(k)+1}$ . This property is illustrated in the following example.

**Example 12.8** Let

$$\mathbf{A} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}.$$

In this case,  $\mathbf{x}^* = [5, 3]^\top$ . Figure 12.4 shows a few iterations of Kaczmarz's algorithm with  $\mu = 1$  and  $\mathbf{x}^{(0)} = \mathbf{0}$ . We have  $\mathbf{a}_1^\top = [1, -1]$ ,  $\mathbf{a}_2^\top = [0, 1]$ ,  $b_1 = 2$ ,  $b_2 = 3$ . In Figure 12.4, the diagonal line passing through the point  $[2, 0]^\top$  corresponds to the set  $\{\mathbf{x} : \mathbf{a}_1^\top \mathbf{x} = b_1\}$ , and the horizontal line passing through the point  $[0, 3]^\top$  corresponds to the set  $\{\mathbf{x} : \mathbf{a}_2^\top \mathbf{x} = b_2\}$ . To illustrate



**Figure 12.4** Iterations of Kaczmarz's algorithm in Example 12.8.

the algorithm, we perform three iterations:

$$\begin{aligned} \mathbf{x}^{(1)} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + (2 - 0) \frac{1}{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \\ \mathbf{x}^{(2)} &= \begin{bmatrix} 1 \\ -1 \end{bmatrix} + (3 - (-1)) \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \\ \mathbf{x}^{(3)} &= \begin{bmatrix} 1 \\ 3 \end{bmatrix} + (2 - (-2)) \frac{1}{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}. \end{aligned}$$

As Figure 12.4 illustrates, the property

$$\mathbf{a}_{R(k)+1}^\top (\mathbf{x}^{(k+1)} - \mathbf{x}^*) = 0$$

holds at every iteration. Note the convergence of the iterations of the algorithm to the solution. ■

## 12.5 Solving Linear Equations in General

Consider the general problem of solving a system of linear equations

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and  $\text{rank } \mathbf{A} = r$ . Note that we always have  $r \leq \min\{m, n\}$ . In the case where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\text{rank } \mathbf{A} = n$ , the unique solution to the equation above has the form  $\mathbf{x}^* = \mathbf{A}^{-1}\mathbf{b}$ . Thus, to solve the problem in this case it is enough to know the inverse  $\mathbf{A}^{-1}$ . In this section we analyze a general approach to solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . The approach involves defining a *pseudoinverse* or *generalized inverse* of a given matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , which plays the role of

$\mathbf{A}^{-1}$  when  $\mathbf{A}$  does not have an inverse (e.g., when  $\mathbf{A}$  is not a square matrix). In particular, we discuss the *Moore-Penrose inverse* of a given matrix  $\mathbf{A}$ , denoted  $\mathbf{A}^\dagger$ .

In our discussion of the Moore-Penrose inverse we use the fact that a nonzero matrix of rank  $r$  can be expressed as the product of a matrix of full column rank  $r$  and a matrix of full row rank  $r$ . Such a factorization is referred to as a *full-rank factorization*, a term which in this context was proposed by Gantmacher [45] and Ben-Israel and Greville [6]. We state and prove the above result in the following lemma.

**Lemma 12.3 Full-Rank Factorization.** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\text{rank } \mathbf{A} = r \leq \min\{m, n\}$ . Then, there exist matrices  $\mathbf{B} \in \mathbb{R}^{m \times r}$  and  $\mathbf{C} \in \mathbb{R}^{r \times n}$  such that*

$$\mathbf{A} = \mathbf{BC},$$

where

$$\text{rank } \mathbf{A} = \text{rank } \mathbf{B} = \text{rank } \mathbf{C} = r.$$

□

*Proof.* Because  $\text{rank } \mathbf{A} = r$ , we can find  $r$  linearly independent columns of  $\mathbf{A}$ . Without loss of generality, let  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$  be such columns, where  $\mathbf{a}_i$  is the  $i$ th column of  $\mathbf{A}$ . The remaining columns of  $\mathbf{A}$  can be expressed as linear combinations of  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r$ . Thus, a possible choice for the full-rank matrices  $\mathbf{B}$  and  $\mathbf{C}$  are

$$\mathbf{B} = [\mathbf{a}_1, \dots, \mathbf{a}_r] \in \mathbb{R}^{m \times r},$$

$$\mathbf{C} = \begin{bmatrix} 1 & \cdots & 0 & c_{1,r+1} & \cdots & c_{1,n} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & c_{r,r+1} & \cdots & c_{r,n} \end{bmatrix} \in \mathbb{R}^{r \times n},$$

where the entries  $c_{i,j}$  are such that for each  $j = r + 1, \dots, n$ , we have  $\mathbf{a}_j = c_{1,j}\mathbf{a}_1 + \cdots + c_{r,j}\mathbf{a}_r$ . Thus,  $\mathbf{A} = \mathbf{BC}$ . ■

Note that if  $m < n$  and  $\text{rank } \mathbf{A} = m$ , then we can take

$$\mathbf{B} = \mathbf{I}_m, \quad \mathbf{C} = \mathbf{A},$$

where  $\mathbf{I}_m$  is the  $m \times m$  identity matrix. If, on the other hand,  $m > n$  and  $\text{rank } \mathbf{A} = n$ , then we can take

$$\mathbf{B} = \mathbf{A}, \quad \mathbf{C} = \mathbf{I}_n.$$

**Example 12.9** Let  $\mathbf{A}$  be given by

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & -2 & 5 \\ 1 & 0 & -3 & 2 \\ 3 & -1 & -13 & 5 \end{bmatrix}.$$

Note that  $\text{rank } \mathbf{A} = 2$ . We can write a full-rank factorization of  $\mathbf{A}$  based on the proof of Lemma 12.3:

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 0 \\ 3 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -3 & 2 \\ 0 & 1 & 4 & 1 \end{bmatrix} = \mathbf{BC}.$$

■

We now introduce the *Moore-Penrose inverse* and discuss its existence and uniqueness. For this, we first consider the matrix equation

$$\mathbf{AXA} = \mathbf{A},$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a given matrix and  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is a matrix we wish to determine. Observe that if  $\mathbf{A}$  is a nonsingular square matrix, then the equation above has the unique solution  $\mathbf{X} = \mathbf{A}^{-1}$ . We now define the Moore-Penrose inverse, also called the pseudoinverse or generalized inverse.

**Definition 12.1** Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , a matrix  $\mathbf{A}^\dagger \in \mathbb{R}^{n \times m}$  is called a *pseudoinverse* of the matrix  $\mathbf{A}$  if

$$\mathbf{AA}^\dagger \mathbf{A} = \mathbf{A},$$

and there exist matrices  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{V} \in \mathbb{R}^{m \times m}$  such that

$$\mathbf{A}^\dagger = \mathbf{U}\mathbf{A}^\top \quad \text{and} \quad \mathbf{A}^\dagger = \mathbf{A}^\top\mathbf{V}.$$

■

The requirement  $\mathbf{A}^\dagger = \mathbf{U}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{V}$  can be interpreted as follows. Each row of the pseudoinverse matrix  $\mathbf{A}^\dagger$  of  $\mathbf{A}$  is a linear combination of the rows of  $\mathbf{A}^\top$ , and each column of  $\mathbf{A}^\dagger$  is a linear combination of the columns of  $\mathbf{A}^\top$ .

For the case in which a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m \geq n$  and  $\text{rank } \mathbf{A} = n$ , we can easily check that the following is a pseudoinverse of  $\mathbf{A}$ :

$$\mathbf{A}^\dagger = (\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top.$$

Indeed,  $\mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top\mathbf{A} = \mathbf{A}$ , and if we define  $\mathbf{U} = (\mathbf{A}^\top\mathbf{A})^{-1}$  and  $\mathbf{V} = \mathbf{A}(\mathbf{A}^\top\mathbf{A})^{-1}(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top$ , then  $\mathbf{A}^\dagger = \mathbf{U}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{V}$ . Note that, in fact, we have  $\mathbf{A}^\dagger\mathbf{A} = \mathbf{I}_n$ . For this reason,  $(\mathbf{A}^\top\mathbf{A})^{-1}\mathbf{A}^\top$  is often called the *left pseudoinverse* of  $\mathbf{A}$ . This formula also appears in least-squares analysis (see Section 12.1).

For the case in which a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $m \leq n$  and  $\text{rank } \mathbf{A} = m$ , we can easily check, as we did in the previous case, that the following is a pseudoinverse of  $\mathbf{A}$ :

$$\mathbf{A}^\dagger = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}.$$

Note that in this case we have  $\mathbf{A}\mathbf{A}^\dagger = \mathbf{I}_m$ . For this reason,  $\mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}$  is often called the *right pseudoinverse* of  $\mathbf{A}$ . This formula also appears in the problem of minimizing  $\|\mathbf{x}\|$  subject to  $\mathbf{Ax} = \mathbf{b}$  (see Section 12.3).

**Theorem 12.4** *Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . If a pseudoinverse  $\mathbf{A}^\dagger$  of  $\mathbf{A}$  exists, then it is unique.*  $\square$

*Proof.* Let  $\mathbf{A}_1^\dagger$  and  $\mathbf{A}_2^\dagger$  be pseudoinverses of  $\mathbf{A}$ . We show that  $\mathbf{A}_1^\dagger = \mathbf{A}_2^\dagger$ . By definition,

$$\mathbf{A}\mathbf{A}_1^\dagger\mathbf{A} = \mathbf{A}\mathbf{A}_2^\dagger\mathbf{A} = \mathbf{A},$$

and there are matrices  $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{n \times n}$ ,  $\mathbf{V}_1, \mathbf{V}_2 \in \mathbb{R}^{m \times m}$  such that

$$\begin{aligned}\mathbf{A}_1^\dagger &= \mathbf{U}_1\mathbf{A}^\top = \mathbf{A}^\top\mathbf{V}_1, \\ \mathbf{A}_2^\dagger &= \mathbf{U}_2\mathbf{A}^\top = \mathbf{A}^\top\mathbf{V}_2.\end{aligned}$$

Let

$$\mathbf{D} = \mathbf{A}_2^\dagger - \mathbf{A}_1^\dagger, \mathbf{U} = \mathbf{U}_2 - \mathbf{U}_1, \mathbf{V} = \mathbf{V}_2 - \mathbf{V}_1.$$

Then, we have

$$\mathbf{O} = \mathbf{A}\mathbf{D}\mathbf{A}, \mathbf{D} = \mathbf{U}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{V}.$$

Therefore, using the two equations above, we have

$$(\mathbf{D}\mathbf{A})^\top\mathbf{D}\mathbf{A} = \mathbf{A}^\top\mathbf{D}^\top\mathbf{D}\mathbf{A} = \mathbf{A}^\top\mathbf{V}^\top\mathbf{A}\mathbf{D}\mathbf{A} = \mathbf{O},$$

which implies that

$$\mathbf{D}\mathbf{A} = \mathbf{O}.$$

On the other hand, because  $\mathbf{D}\mathbf{A} = \mathbf{O}$ , we have

$$\mathbf{D}\mathbf{D}^\top = \mathbf{D}\mathbf{A}\mathbf{U}^\top = \mathbf{O},$$

which implies that

$$\mathbf{D} = \mathbf{A}_2^\dagger - \mathbf{A}_1^\dagger = \mathbf{O}$$

and hence

$$\mathbf{A}_2^\dagger = \mathbf{A}_1^\dagger.$$

■

From Theorem 12.4, we know that if a pseudoinverse matrix exists, then it is unique. Our goal now is to show that the pseudoinverse always exists. In fact, we show that the pseudoinverse of any given matrix  $\mathbf{A}$  is given by the formula

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger\mathbf{B}^\dagger,$$

where  $\mathbf{B}^\dagger$  and  $\mathbf{C}^\dagger$  are the pseudoinverses of the matrices  $\mathbf{B}$  and  $\mathbf{C}$  that form a full-rank factorization of  $\mathbf{A}$ ; that is,  $\mathbf{A} = \mathbf{BC}$ , where  $\mathbf{B}$  and  $\mathbf{C}$  are of full

rank (see Lemma 12.3). Note that we already know how to compute  $\mathbf{B}^\dagger$  and  $\mathbf{C}^\dagger$ :

$$\mathbf{B}^\dagger = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top, \quad \mathbf{C}^\dagger = \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1}.$$

**Theorem 12.5** *Let a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  have a full-rank factorization  $\mathbf{A} = \mathbf{BC}$ , with  $\text{rank } \mathbf{A} = \text{rank } \mathbf{B} = \text{rank } \mathbf{C} = r$ ,  $\mathbf{B} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{C} \in \mathbb{R}^{r \times n}$ . Then,*

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger.$$

□

*Proof.* We show that

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger = \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$$

satisfies the conditions of Definition 12.1 for a pseudoinverse. Indeed, first observe that

$$\mathbf{AC}^\dagger \mathbf{B}^\dagger \mathbf{A} = \mathbf{BCC}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{BC} = \mathbf{BC} = \mathbf{A}.$$

Next, define

$$\mathbf{U} = \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} (\mathbf{C} \mathbf{C}^\top)^{-1} \mathbf{C}$$

and

$$\mathbf{V} = \mathbf{B} (\mathbf{B}^\top \mathbf{B})^{-1} (\mathbf{C} \mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top.$$

It is easy to verify that the matrices  $\mathbf{U}$  and  $\mathbf{V}$  above satisfy

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger = \mathbf{U} \mathbf{A}^\top = \mathbf{A}^\top \mathbf{V}.$$

Hence,

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger$$

is the pseudoinverse of  $\mathbf{A}$ . ■

**Example 12.10** *Continued from Example 12.9.* Recall that

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & -2 & 5 \\ 1 & 0 & -3 & 2 \\ 3 & -1 & -13 & 5 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 0 \\ 3 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -3 & 2 \\ 0 & 1 & 4 & 1 \end{bmatrix} = \mathbf{BC}.$$

We compute

$$\mathbf{B}^\dagger = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top = \frac{1}{27} \begin{bmatrix} 5 & 2 & 5 \\ 16 & 1 & -11 \end{bmatrix}$$

and

$$\mathbf{C}^\dagger = \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} = \frac{1}{76} \begin{bmatrix} 9 & 5 \\ 5 & 7 \\ -7 & 13 \\ 23 & 17 \end{bmatrix}.$$

Thus,

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger = \frac{1}{2052} \begin{bmatrix} 125 & 23 & -10 \\ 137 & 17 & -52 \\ 173 & -1 & -178 \\ 387 & 63 & -72 \end{bmatrix}.$$

■

We emphasize that the formula  $\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger$  does not necessarily hold if  $\mathbf{A} = \mathbf{BC}$  is not a full-rank factorization. The following example, which is taken from [45], illustrates this point.

**Example 12.11** Let

$$\mathbf{A} = \begin{bmatrix} 1 \end{bmatrix}.$$

Obviously,  $\mathbf{A}^\dagger = \mathbf{A}^{-1} = \mathbf{A} = [1]$ . Observe that  $\mathbf{A}$  can be represented as

$$\mathbf{A} = [0 \ 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{BC}.$$

The above is not a full-rank factorization of  $\mathbf{A}$ . Let us now compute  $\mathbf{B}^\dagger$  and  $\mathbf{C}^\dagger$ . We have

$$\mathbf{B}^\dagger = \mathbf{B}^\top (\mathbf{BB}^\top)^{-1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

$$\mathbf{C}^\dagger = (\mathbf{C}^\top \mathbf{C})^{-1} \mathbf{C}^\top = \begin{bmatrix} 1/2 & 1/2 \end{bmatrix}.$$

(Note that the formulas for  $\mathbf{B}^\dagger$  and  $\mathbf{C}^\dagger$  here are different from those in Example 12.10 because of the dimensions of  $\mathbf{B}$  and  $\mathbf{C}$  in this example.) Thus,

$$\mathbf{C}^\dagger \mathbf{B}^\dagger = \begin{bmatrix} 1/2 \end{bmatrix},$$

which is not equal to  $\mathbf{A}^\dagger$ .

■

We can simplify the expression

$$\mathbf{A}^\dagger = \mathbf{C}^\dagger \mathbf{B}^\dagger = \mathbf{C}^\top (\mathbf{CC}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$$

to

$$\mathbf{A}^\dagger = \mathbf{C}^\top (\mathbf{B}^\top \mathbf{AC}^\top)^{-1} \mathbf{B}^\top.$$

The expression above is easily verified simply by substituting  $\mathbf{A} = \mathbf{BC}$ . This explicit formula for  $\mathbf{A}^\dagger$  is attributed to C. C. MacDuffee by Ben-Israel and Greville [6]. Ben-Israel and Greville report that around 1959, MacDuffee was the first to point out that a full-rank factorization of  $\mathbf{A}$  leads to the above

explicit formula. However, they mention that MacDuffee did it in a private communication, so there is no published work by MacDuffee that contains the result.

We now prove two important properties of  $\mathbf{A}^\dagger$  in the context of solving a system of linear equations  $\mathbf{Ax} = \mathbf{b}$ .

**Theorem 12.6** Consider a system of linear equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\text{rank } \mathbf{A} = r$ . The vector  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b}$  minimizes  $\|\mathbf{Ax} - \mathbf{b}\|^2$  on  $\mathbb{R}^n$ . Furthermore, among all vectors in  $\mathbb{R}^n$  that minimize  $\|\mathbf{Ax} - \mathbf{b}\|^2$ , the vector  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b}$  is the unique vector with minimal norm.  $\square$

*Proof.* We first show that  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b}$  minimizes  $\|\mathbf{Ax} - \mathbf{b}\|^2$  over  $\mathbb{R}^n$ . To this end, observe that for any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\begin{aligned}\|\mathbf{Ax} - \mathbf{b}\|^2 &= \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*) + \mathbf{Ax}^* - \mathbf{b}\|^2 \\ &= \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 + \|\mathbf{Ax}^* - \mathbf{b}\|^2 + 2[\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{Ax}^* - \mathbf{b}).\end{aligned}$$

We now show that

$$[\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{Ax}^* - \mathbf{b}) = 0.$$

Indeed,

$$\begin{aligned}[\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{Ax}^* - \mathbf{b}) &= (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{A}^\top \mathbf{Ax}^* - \mathbf{A}^\top \mathbf{b}) \\ &= (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{A}^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{b} - \mathbf{A}^\top \mathbf{b}).\end{aligned}$$

However,

$$\mathbf{A}^\top \mathbf{A} \mathbf{A}^\dagger = \mathbf{C}^\top \mathbf{B}^\top \mathbf{B} \mathbf{C} \mathbf{C}^\top (\mathbf{C} \mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top = \mathbf{A}^\top.$$

Hence,

$$[\mathbf{A}(\mathbf{x} - \mathbf{x}^*)]^\top (\mathbf{Ax}^* - \mathbf{b}) = (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{A}^\top \mathbf{b} - \mathbf{A}^\top \mathbf{b}) = 0.$$

Thus, we have

$$\|\mathbf{Ax} - \mathbf{b}\|^2 = \|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 + \|\mathbf{Ax}^* - \mathbf{b}\|^2.$$

Because

$$\|\mathbf{A}(\mathbf{x} - \mathbf{x}^*)\|^2 \geq 0,$$

we obtain

$$\|\mathbf{Ax} - \mathbf{b}\|^2 \geq \|\mathbf{Ax}^* - \mathbf{b}\|^2$$

and thus  $\mathbf{x}^*$  minimizes  $\|\mathbf{Ax} - \mathbf{b}\|^2$ .

We now show that among all  $\mathbf{x}$  that minimize  $\|\mathbf{Ax} - \mathbf{b}\|^2$ , the vector  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b}$  is the unique vector with minimum norm. So let  $\tilde{\mathbf{x}}$  be a vector minimizing  $\|\mathbf{Ax} - \mathbf{b}\|^2$ . We have

$$\begin{aligned}\|\tilde{\mathbf{x}}\|^2 &= \|(\tilde{\mathbf{x}} - \mathbf{x}^*) + \mathbf{x}^*\|^2 \\ &= \|\tilde{\mathbf{x}} - \mathbf{x}^*\|^2 + \|\mathbf{x}^*\|^2 + 2\mathbf{x}^{*\top}(\tilde{\mathbf{x}} - \mathbf{x}^*).\end{aligned}$$

Observe that

$$\mathbf{x}^{*\top}(\tilde{\mathbf{x}} - \mathbf{x}^*) = 0.$$

To see this, note that

$$\begin{aligned}\mathbf{x}^{*\top}(\tilde{\mathbf{x}} - \mathbf{x}^*) &= (\mathbf{A}^\dagger \mathbf{b})^\top(\tilde{\mathbf{x}} - \mathbf{A}^\dagger \mathbf{b}) \\ &= \mathbf{b}^\top \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-\top}(\mathbf{C} \mathbf{C}^\top)^{-\top} \mathbf{C}(\tilde{\mathbf{x}} - \mathbf{C}^\top(\mathbf{C} \mathbf{C}^\top)^{-1}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{b}) \\ &= \mathbf{b}^\top \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-\top}(\mathbf{C} \mathbf{C}^\top)^{-\top}[\mathbf{C} \tilde{\mathbf{x}} - (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{b}],\end{aligned}$$

where the superscript  $-\top$  denotes the transpose of the inverse. Now,  $\|\mathbf{Ax} - \mathbf{b}\|^2 = \|\mathbf{B}(\mathbf{Cx}) - \mathbf{b}\|^2$ . Because  $\tilde{\mathbf{x}}$  minimizes  $\|\mathbf{Ax} - \mathbf{b}\|^2$  and  $\mathbf{C}$  is of full rank, then  $\mathbf{y}^* = \mathbf{C}\tilde{\mathbf{x}}$  minimizes  $\|\mathbf{By} - \mathbf{b}\|^2$  over  $\mathbb{R}^r$  (see Exercise 12.29). Because  $\mathbf{B}$  is of full rank, by Theorem 12.1, we have  $\mathbf{C}\tilde{\mathbf{x}} = \mathbf{y}^* = (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{b}$ . Substituting this into the equation above, we get  $\mathbf{x}^{*\top}(\tilde{\mathbf{x}} - \mathbf{x}^*) = 0$ .

Therefore, we have

$$\|\tilde{\mathbf{x}}\|^2 = \|\mathbf{x}^*\|^2 + \|\tilde{\mathbf{x}} - \mathbf{x}^*\|^2.$$

For all  $\tilde{\mathbf{x}} \neq \mathbf{x}^*$ , we have

$$\|\tilde{\mathbf{x}} - \mathbf{x}^*\|^2 > 0,$$

and hence

$$\|\tilde{\mathbf{x}}\|^2 > \|\mathbf{x}^*\|^2$$

or, equivalently,

$$\|\tilde{\mathbf{x}}\| > \|\mathbf{x}^*\|.$$

Hence, among all vectors minimizing  $\|\mathbf{Ax} - \mathbf{b}\|^2$ , the vector  $\mathbf{x}^* = \mathbf{A}^\dagger \mathbf{b}$  is the unique vector with minimum norm. ■

The generalized inverse has the following useful properties (see Exercise 12.30):

- a.  $(\mathbf{A}^\top)^\dagger = (\mathbf{A}^\dagger)^\top$ .
- b.  $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$ .

These two properties are similar to those that are satisfied by the usual matrix inverse. However, we point out that the property  $(\mathbf{A}_1 \mathbf{A}_2)^\dagger = \mathbf{A}_2^\dagger \mathbf{A}_1^\dagger$  does not hold in general (see Exercise 12.32).

Finally, it is important to note that we can define the generalized inverse in an alternative way, following the definition of Penrose. Specifically, the Penrose definition of the generalized inverse of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is the unique matrix  $\mathbf{A}^\dagger \in \mathbb{R}^{n \times m}$  satisfying the following properties:

1.  $\mathbf{A}\mathbf{A}^\dagger \mathbf{A} = \mathbf{A}$ .
2.  $\mathbf{A}^\dagger \mathbf{A}\mathbf{A}^\dagger = \mathbf{A}^\dagger$ .

$$3. (\mathbf{A}\mathbf{A}^\dagger)^\top = \mathbf{A}\mathbf{A}^\dagger.$$

$$4. (\mathbf{A}^\dagger\mathbf{A})^\top = \mathbf{A}^\dagger\mathbf{A}.$$

The Penrose definition above is equivalent to Definition 12.1 (see Exercise 12.31). For more information on generalized inverses and their applications, we refer the reader to the books by Ben-Israel and Greville [6] and Campbell and Meyer [23].

## EXERCISES

**12.1** A rock is accelerated to 3, 5, and 6 m/s<sup>2</sup> by applying forces of 1, 2, and 3 N, respectively. Assuming that Newton's law  $F = ma$  holds, where  $F$  is the force and  $a$  is the acceleration, estimate the mass  $m$  of the rock using the least-squares method.

**12.2** A spring is stretched to lengths  $L = 3, 4$ , and 5 cm under applied forces  $F = 1, 2$ , and 4 N, respectively. Assuming that Hooke's law  $L = a + bF$  holds, estimate the normal length  $a$  and spring constant  $b$  using the least-squares approach.

**12.3** Suppose that we perform an experiment to calculate the gravitational constant  $g$  as follows. We drop a ball from a certain height and measure its distance from the original point at certain time instants. The results of the experiment are shown in the following table.

Time (seconds)	1.00	2.00	3.00
Distance (meters)	5.00	19.5	44.0

The equation relating the distance  $s$  and the time  $t$  at which  $s$  is measured is given by

$$s = \frac{1}{2}gt^2.$$

- a. Find a least-squares estimate of  $g$  using the experimental results from the table above.
- b. Suppose that we take an additional measurement at time 4.00 and obtain a distance of 78.5. Use the recursive least-squares algorithm to calculate an updated least-squares estimate of  $g$ .

**12.4** Suppose that we have a speech signal, represented as a finite sequence of real numbers  $x_1, x_2, \dots, x_n$ . Suppose that we record this signal onto magnetic tape. The recorded speech signal is represented by another sequence of real numbers  $y_1, y_2, \dots, y_n$ .

Suppose that we model the recording process as a simple scaling of the original signal (i.e., we believe that a good model of the relationship between the recorded signal and the original signal is  $y_i = \alpha x_i$  for some constant  $\alpha$  that does not depend on  $i$ ). Suppose that we know exactly the original signal  $x_1, x_2, \dots, x_n$  as well as the recorded signal  $y_1, y_2, \dots, y_n$ . Use the least-squares method to find a formula for estimating the scale factor  $\alpha$  given this data. (You may assume that at least one  $x_i$  is nonzero.)

**12.5** Suppose that we wish to estimate the value of the resistance  $R$  of a resistor. Ohm's law states that if  $V$  is the voltage across the resistor and  $I$  is the current through the resistor, then  $V = IR$ . To estimate  $R$ , we apply a 1-ampere current through the resistor and measure the voltage across it. We perform the experiment on  $n$  voltage-measuring devices and obtain measurements of  $V_1, \dots, V_n$ . Show that the least-squares estimate of  $R$  is simply the average of  $V_1, \dots, V_n$ .

**12.6** The table below shows the stock prices for three companies, X, Y, and Z, tabulated over three days:

	Day		
	1	2	3
X	6	4	5
Y	1	1	3
Z	2	1	2

Suppose that an investment analyst proposes a model for the predicting the stock price of X based on those of Y and Z:

$$p_X = ap_Y + bp_Z,$$

where  $p_X$ ,  $p_Y$ , and  $p_Z$  are the stock prices of X, Y, and Z, respectively, and  $a$  and  $b$  are real-valued parameters. Calculate the least-squares estimate of parameters  $a$  and  $b$  based on the data in the table above.

**12.7** We are given two mixtures, A and B. Mixture A contains 30% gold, 40% silver, and 30% platinum, whereas mixture B contains 10% gold, 20% silver, and 70% platinum (all percentages of weight). We wish to determine the ratio of the weight of mixture A to the weight of mixture B such that we have as close as possible to a total of 5 ounces of gold, 3 ounces of silver, and 4 ounces of platinum. Formulate and solve the problem using the linear least-squares method.

**12.8** *Background:* If  $\mathbf{A}\mathbf{x} + \mathbf{w} = \mathbf{b}$ , where  $\mathbf{w}$  is a “white noise” vector, then define the *least-squares estimate* of  $\mathbf{x}$  given  $\mathbf{b}$  to be the solution to the problem

$$\text{minimize } \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2.$$

This problem is related to Wiener filtering.

*Application:* Suppose that a given speech signal  $\{u_k : k = 1, \dots, n\}$  (with  $u_k \in \mathbb{R}$ ) is transmitted over a telephone cable with input-output behavior given by  $y_k = ay_{k-1} + bu_k + v_k$ , where, at each time  $k$ ,  $y_k \in \mathbb{R}$  is the output,  $u_k \in \mathbb{R}$  is the input (speech signal value), and  $v_k$  represents white noise. The parameters  $a$  and  $b$  are fixed known constants, and the initial condition is  $y_0 = 0$ .

We can measure the signal  $\{y_k\}$  at the output of the telephone cable, but we cannot directly measure the desired signal  $\{u_k\}$  or the noise signal  $\{v_k\}$ . Derive a formula for the linear least-squares estimate of the signal  $\{u_k : k = 1, \dots, n\}$  given the signal  $\{y_k : k = 1, \dots, n\}$ .

*Note:* Even though the vector  $\mathbf{v} = [v_1, \dots, v_n]^\top$  is a white noise vector, the vector  $\mathbf{D}\mathbf{v}$  (where  $\mathbf{D}$  is a matrix) is, in general, not.

**12.9 Line Fitting.** Let  $[x_1, y_1]^\top, \dots, [x_p, y_p]^\top$ ,  $p \geq 2$ , be points in  $\mathbb{R}^2$ . We wish to find the straight line of best fit through these points (“best” in the sense that the total squared error is minimized); that is, we wish to find  $a^*, b^* \in \mathbb{R}$  to minimize

$$f(a, b) = \sum_{i=1}^p (ax_i + b - y_i)^2.$$

Assume that the  $x_i$ ,  $i = 1, \dots, p$ , are not all equal. Show that there exist unique parameters  $a^*$  and  $b^*$  for the line of best fit, and find the parameters in terms of the following quantities:

$$\begin{aligned}\overline{X} &= \frac{1}{p} \sum_{i=1}^p x_i, \\ \overline{Y} &= \frac{1}{p} \sum_{i=1}^p y_i, \\ \overline{X^2} &= \frac{1}{p} \sum_{i=1}^p x_i^2, \\ \overline{Y^2} &= \frac{1}{p} \sum_{i=1}^p y_i^2, \\ \overline{XY} &= \frac{1}{p} \sum_{i=1}^p x_i y_i.\end{aligned}$$

**12.10** Suppose that we take measurements of a sinusoidal signal  $y(t) = \sin(\omega t + \theta)$  at times  $t_1, \dots, t_p$ , and obtain values  $y_1, \dots, y_p$ , where  $-\pi/2 \leq \omega t_i + \theta \leq \pi/2$ ,  $i = 1, \dots, p$ , and the  $t_i$  are not all equal. We wish to determine the values of the frequency  $\omega$  and phase  $\theta$ .

- a. Express the problem as a system of linear equations.
- b. Find the least-squares estimate of  $\omega$  and  $\theta$  based on part a. Use the following notation:

$$\begin{aligned}\bar{T} &= \frac{1}{p} \sum_{i=1}^p t_i, \\ \bar{T^2} &= \frac{1}{p} \sum_{i=1}^p t_i^2, \\ \bar{TY} &= \frac{1}{p} \sum_{i=1}^p t_i \arcsin y_i, \\ \bar{Y} &= \frac{1}{p} \sum_{i=1}^p \arcsin y_i.\end{aligned}$$

**12.11** We are given a point  $[x_0, y_0]^\top \in \mathbb{R}^2$ . Consider the straight line on the  $\mathbb{R}^2$  plane given by the equation  $y = mx$ . Using a least-squares formulation, find the point on the straight line that is closest to the given point  $[x_0, y_0]$ , where the measure of closeness is in terms of the Euclidean norm on  $\mathbb{R}^2$ .

*Hint:* The given line can be expressed as the range of the matrix  $\mathbf{A} = [1, m]^\top$ .

**12.12** Consider the affine function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  of the form  $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + c$ , where  $\mathbf{a} \in \mathbb{R}^n$  and  $c \in \mathbb{R}$ .

- a. We are given a set of  $p$  pairs  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_p, y_p)$ , where  $\mathbf{x}_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, \dots, p$ . We wish to find the affine function of best fit to these points, where “best” is in the sense of minimizing the total square error

$$\sum_{i=1}^p (f(\mathbf{x}_i) - y_i)^2.$$

Formulate the above as an optimization problem of the form: minimize  $\|\mathbf{Az} - \mathbf{b}\|^2$  with respect to  $\mathbf{z}$ . Specify the dimensions of  $\mathbf{A}$ ,  $\mathbf{z}$ , and  $\mathbf{b}$ .

- b. Suppose that the points satisfy

$$\mathbf{x}_1 + \cdots + \mathbf{x}_p = \mathbf{0}$$

and

$$y_1 \mathbf{x}_1 + \cdots + y_p \mathbf{x}_p = \mathbf{0}.$$

Find the affine function of best fit in this case, assuming that it exists and is unique.



**Figure 12.5** Input-output system in Exercise 12.13.

**12.13** For the system shown in Figure 12.5, consider a set of input-output pairs  $(u_1, y_1), \dots, (u_n, y_n)$ , where  $u_k \in \mathbb{R}$ ,  $y_k \in \mathbb{R}$ ,  $k = 1, \dots, n$ .

- Suppose that we wish to find the best linear estimate of the system based on the input-output data above. In other words, we wish to find a  $\hat{\theta}_n \in \mathbb{R}$  to fit the model  $y_k = \hat{\theta}_n u_k$ ,  $k = 1, \dots, n$ . Using the least-squares approach, derive a formula for  $\hat{\theta}_n$  based on  $u_1, \dots, u_n$  and  $y_1, \dots, y_n$ .
- Suppose that the data in part a are generated by

$$y_k = \theta u_k + e_k,$$

where  $\theta \in \mathbb{R}$  and  $u_k = 1$  for all  $k$ . Show that the parameter  $\hat{\theta}_n$  in part a converges to  $\theta$  as  $n \rightarrow \infty$  if and only if

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n e_k = 0.$$

**12.14** Consider a discrete-time linear system  $x_{k+1} = ax_k + bu_k$ , where  $u_k$  is the input at time  $k$ ,  $x_k$  is the output at time  $k$ , and  $a, b \in \mathbb{R}$  are system parameters. Suppose that we apply a constant input  $u_k = 1$  for all  $k \geq 0$  and measure the first four values of the output to be  $x_0 = 0$ ,  $x_1 = 1$ ,  $x_2 = 2$ ,  $x_3 = 8$ . Find the least-squares estimate of  $a$  and  $b$  based on the data above.

**12.15** Consider a discrete-time linear system  $x_{k+1} = ax_k + bu_k$ , where  $u_k$  is the input at time  $k$ ,  $x_k$  is the output at time  $k$ , and  $a, b \in \mathbb{R}$  are system parameters. Given the first  $n + 1$  values of the impulse response  $h_0, \dots, h_n$ , find the least-squares estimate of  $a$  and  $b$ . You may assume that at least one  $h_k$  is nonzero.

*Note:* The *impulse response* is the output sequence resulting from an input of  $u_0 = 1$ ,  $u_k = 0$  for  $k \neq 0$  and zero initial condition  $x_0 = 0$ .

**12.16** Consider a discrete-time linear system  $x_{k+1} = ax_k + bu_k$ , where  $u_k$  is the input at time  $k$ ,  $x_k$  is the output at time  $k$ , and  $a, b \in \mathbb{R}$  are system parameters. Given the first  $n + 1$  values of the step response  $s_0, \dots, s_n$ , where  $n > 1$ , find the least-squares estimate of  $a$  and  $b$ . You may assume that at least one  $s_k$  is nonzero.

*Note:* The *step response* is the output sequence resulting from an input of  $u_k = 1$  for  $k \geq 0$ , and zero initial condition  $x_0 = 0$  (i.e.,  $s_0 = x_0 = 0$ ).

**12.17** Consider a known discrete-time signal on the time interval  $\{1, \dots, n\}$ , represented by the vector  $\mathbf{x} \in \mathbb{R}^n$  ( $x_i$  is the value of the signal at time  $i$ ). We transmit the signal  $a\mathbf{x}$  over a communication channel, where  $a \in \mathbb{R}$  represents the “amplitude” of the transmission, a quantity unknown to the receiver. The receiver receives a signal  $\mathbf{y} \in \mathbb{R}^n$ , which is a distorted version of the transmitted signal (so that  $\mathbf{y}$  may not be equal to  $a\mathbf{x}$  for any  $a$ ). Formulate the problem of estimating the quantity  $a$  according to a least-squares criterion, and solve it (stating whatever appropriate assumptions are necessary, if any).

**12.18** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \geq n$ , and  $\text{rank } \mathbf{A} = n$ . Consider the constrained optimization problem

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top \mathbf{b} \\ & \text{subject to} \quad \mathbf{x} \in \mathcal{R}(\mathbf{A}), \end{aligned}$$

where  $\mathcal{R}(\mathbf{A})$  denotes the range of  $\mathbf{A}$ . Derive an expression for the global minimizer of this problem in terms of  $\mathbf{A}$  and  $\mathbf{b}$ .

**12.19** Solve the problem

$$\begin{aligned} & \text{minimize} \quad \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} \quad \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \mathbf{x} = 1, \end{aligned}$$

where  $\mathbf{x}_0 = [0, -3, 0]^\top$ .

**12.20** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \leq n$ ,  $\text{rank } \mathbf{A} = m$ , and  $\mathbf{x}_0 \in \mathbb{R}^n$ . Consider the problem

$$\begin{aligned} & \text{minimize} \quad \|\mathbf{x} - \mathbf{x}_0\| \\ & \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned}$$

Show that this problem has a unique solution given by

$$\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b} + (\mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}) \mathbf{x}_0.$$

**12.21** Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ ,  $\text{rank } \mathbf{A} = n$ , and  $\mathbf{b}_1, \dots, \mathbf{b}_p \in \mathbb{R}^m$ , consider the problem

$$\text{minimize} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}_1\|^2 + \|\mathbf{A}\mathbf{x} - \mathbf{b}_2\|^2 + \cdots + \|\mathbf{A}\mathbf{x} - \mathbf{b}_p\|^2.$$

Suppose that  $\mathbf{x}_i^*$  is a solution to the problem

$$\text{minimize} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}_i\|^2,$$

where  $i = 1, \dots, p$ . Write the solution to the problem in terms of  $\mathbf{x}_1^*, \dots, \mathbf{x}_p^*$ .

**12.22** Given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ ,  $\text{rank } \mathbf{A} = n$ ,  $\mathbf{b}_1, \dots, \mathbf{b}_p \in \mathbb{R}^m$ , and  $\alpha_1, \dots, \alpha_p \in \mathbb{R}$ , consider the problem

$$\text{minimize } \alpha_1 \|\mathbf{Ax} - \mathbf{b}_1\|^2 + \alpha_2 \|\mathbf{Ax} - \mathbf{b}_2\|^2 + \cdots + \alpha_p \|\mathbf{Ax} - \mathbf{b}_p\|^2.$$

Suppose that  $\mathbf{x}_i^*$  is the solution to the problem

$$\text{minimize } \|\mathbf{Ax} - \mathbf{b}_i\|^2,$$

where  $i = 1, \dots, p$ . Assuming that  $\alpha_1 + \cdots + \alpha_p > 0$ , derive a simple expression for the solution to this problem in terms of  $\mathbf{x}_1^*, \dots, \mathbf{x}_p^*$  and  $\alpha_1, \dots, \alpha_p$ .

**12.23** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . Show that  $\mathbf{x}^* = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{b}$  is the only vector in  $\mathcal{R}(\mathbf{A}^\top)$  satisfying  $\mathbf{Ax}^* = \mathbf{b}$ .

**12.24** The purpose of this question is to derive a recursive least-squares algorithm where we *remove* (instead of add) a data point. To formulate the algorithm, suppose that we are given matrices  $\mathbf{A}_0$  and  $\mathbf{A}_1$  such that

$$\mathbf{A}_0 = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{a}_1^\top \end{bmatrix},$$

where  $\mathbf{a}_1 \in \mathbb{R}^n$ . Similarly, suppose that vectors  $\mathbf{b}^{(0)}$  and  $\mathbf{b}^{(1)}$  satisfy

$$\mathbf{b}^{(0)} = \begin{bmatrix} \mathbf{b}^{(1)} \\ b_1 \end{bmatrix},$$

where  $b_1 \in \mathbb{R}$ . Let  $\mathbf{x}^{(0)}$  be the least-squares solution associated with  $(\mathbf{A}_0, \mathbf{b}^{(0)})$  and  $\mathbf{x}^{(1)}$  the least-squares solution associated with  $(\mathbf{A}_1, \mathbf{b}^{(1)})$ . Our goal is to write  $\mathbf{x}^{(1)}$  in terms of  $\mathbf{x}^{(0)}$  and the data point “removed,”  $(\mathbf{a}_1, b_1)$ . As usual, let  $\mathbf{G}_0$  and  $\mathbf{G}_1$  be the Grammians associated with  $\mathbf{x}^{(0)}$  and  $\mathbf{x}^{(1)}$ , respectively.

- a. Write down expressions for the least-squares solutions  $\mathbf{x}^{(0)}$  and  $\mathbf{x}^{(1)}$  in terms of  $\mathbf{A}_0$ ,  $\mathbf{b}^{(0)}$ ,  $\mathbf{A}_1$ , and  $\mathbf{b}^{(1)}$ .
- b. Derive a formula for  $\mathbf{G}_1$  in terms of  $\mathbf{G}_0$  and  $\mathbf{a}_1$ .
- c. Let  $\mathbf{P}_0 = \mathbf{G}_0^{-1}$  and  $\mathbf{P}_1 = \mathbf{G}_1^{-1}$ . Derive a formula for  $\mathbf{P}_1$  in terms of  $\mathbf{P}_0$  and  $\mathbf{a}_1$ . (The formula must not contain any matrix inversions.)
- d. Derive a formula for  $\mathbf{A}_0^\top \mathbf{b}^{(0)}$  in terms of  $\mathbf{G}_1$ ,  $\mathbf{x}^{(0)}$ , and  $\mathbf{a}_1$ .
- e. Finally, derive a formula for  $\mathbf{x}^{(1)}$  in terms of  $\mathbf{x}^{(0)}$ ,  $\mathbf{P}_1$ ,  $\mathbf{a}_1$ , and  $b_1$ . Use this and part c to write a recursive algorithm associated with successive removals of rows from  $(\mathbf{A}_k, \mathbf{b}^{(k)})$ .

**12.25** Show that in Kaczmarz's algorithm, if  $\mathbf{x}^{(0)} = \mathbf{0}$ , then  $\mathbf{x}^{(k)} \in \mathcal{R}(\mathbf{A}^\top)$  for all  $k$ .

**12.26** Consider Kaczmarz's algorithm with  $\mathbf{x}^{(0)} \neq \mathbf{0}$ .

- a. Show that there exists a unique point minimizing  $\|\mathbf{x} - \mathbf{x}^{(0)}\|$  subject to  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ .
- b. Show that Kaczmarz's algorithm converges to the point in part a.

**12.27** Consider Kaczmarz's algorithm with  $\mathbf{x}^{(0)} = \mathbf{0}$ , where  $m = 1$ ; that is,  $\mathbf{A} = [\mathbf{a}^\top] \in \mathbb{R}^{1 \times n}$  and  $\mathbf{a} \neq \mathbf{0}$ , and  $0 < \mu < 2$ . Show that there exists  $0 \leq \gamma < 1$  such that  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \gamma \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$  for all  $k \geq 0$ .

**12.28** Show that in Kaczmarz's algorithm, if  $\mu = 1$ , then  $b_{R(k)+1} - \mathbf{a}_{R(k)+1}^\top \mathbf{x}^{(k+1)} = 0$  for each  $k$ .

**12.29** Consider the problem of minimizing  $\|\mathbf{Ax} - \mathbf{b}\|^2$  over  $\mathbb{R}^n$ , where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ . Let  $\mathbf{x}^*$  be a solution. Suppose that  $\mathbf{A} = \mathbf{BC}$  is a full-rank factorization of  $\mathbf{A}$ ; that is,  $\text{rank } \mathbf{A} = \text{rank } \mathbf{B} = \text{rank } \mathbf{C} = r$ , and  $\mathbf{B} \in \mathbb{R}^{m \times r}$ ,  $\mathbf{C} \in \mathbb{R}^{r \times n}$ . Show that the minimizer of  $\|\mathbf{By} - \mathbf{b}\|$  over  $\mathbb{R}^r$  is  $\mathbf{Cx}^*$ .

**12.30** Prove the following properties of generalized inverses:

- a.  $(\mathbf{A}^\top)^\dagger = (\mathbf{A}^\dagger)^\top$ .
- b.  $(\mathbf{A}^\dagger)^\dagger = \mathbf{A}$ .

**12.31** Show that the Penrose definition of the generalized inverse is equivalent to Definition 12.1.

**12.32** Construct matrices  $\mathbf{A}_1$  and  $\mathbf{A}_2$  such that  $(\mathbf{A}_1 \mathbf{A}_2)^\dagger \neq \mathbf{A}_2^\dagger \mathbf{A}_1^\dagger$ .



# CHAPTER 13

---

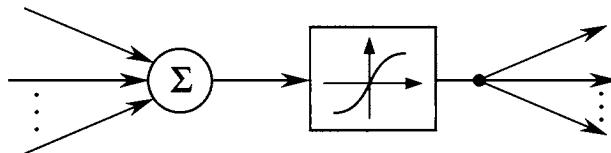
## UNCONSTRAINED OPTIMIZATION AND NEURAL NETWORKS

---

### 13.1 Introduction

In this chapter we apply the techniques of previous chapters to the training of feedforward neural networks. Neural networks have found numerous practical applications, ranging from telephone echo cancellation to aiding in the interpretation of EEG data (see, e.g., [108] and [72]). The essence of neural networks lies in the connection weights between neurons. The selection of these weights is referred to as *training* or *learning*. For this reason, we often refer to the weights as the *learning parameters*. A popular method for training a neural network is the *backpropagation algorithm*, based on an unconstrained optimization problem and an associated gradient algorithm applied to the problem. This chapter is devoted to a description of neural networks and the use of techniques developed in preceding chapters for the training of neural networks.

An *artificial neural network* is a circuit composed of interconnected simple circuit elements called *neurons*. Each neuron represents a map, typically with multiple inputs and a single output. Specifically, the output of the neuron is a function of the sum of the inputs, as illustrated in Figure 13.1. The

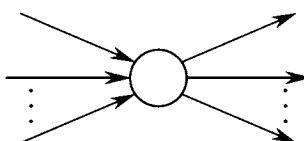


**Figure 13.1** Single neuron.

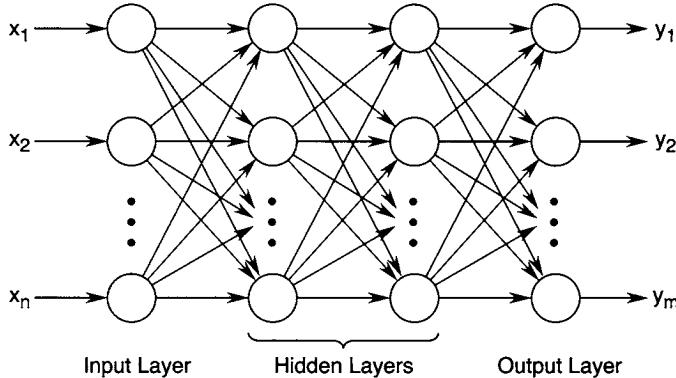
function at the output of the neuron is called the *activation function*. We use the symbol shown in Figure 13.2 to represent a single neuron. Note that the single output of the neuron may be used as an input to several other neurons, and therefore the symbol for a single neuron has multiple arrows emanating from it. A neural network may be implemented using an analog circuit. In this case inputs and outputs may be represented by currents and voltages.

A neural network consists of interconnected neurons, with the inputs to each neuron consisting of weighted outputs of other neurons. The interconnections allow exchange of data or information between neurons. In a *feed-forward neural network*, the neurons are interconnected in layers, so that the data flow in only one direction. Thus, each neuron receives information only from neurons in the preceding layer: The inputs to each neuron are weighted outputs of neurons in the preceding layer. Figure 13.3 illustrates the structure of feedforward neural networks. The first layer in the network is called the *input layer*, and the last layer is called the *output layer*. The layers in between the input and output layers are called *hidden layers*.

We can view a neural network as simply a particular implementation of a map from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ , where  $n$  is the number of inputs  $x_1, \dots, x_n$  and  $m$  is the number of outputs  $y_1, \dots, y_m$ . The map that is implemented by a neural network depends on the weights of the interconnections in the network. Therefore, we can change the mapping that is implemented by the network by adjusting the values of the weights in the network. The information about the mapping is “stored” in the weights over all the neurons, and thus the neural network is a *distributed* representation of the mapping. Moreover, for any given input, computation of the corresponding output is achieved through the collective effect of individual input-output characteristics of each neuron; therefore, the neural network can be considered as a *parallel* computation device. We point out that the ability to implement or approximate a map encompasses many important practical applications. For example, pattern



**Figure 13.2** Symbol for a single neuron.



**Figure 13.3** Structure of a feedforward neural network.

recognition and classification problems can be viewed as function implementation or approximation problems.

Suppose that we are given a map  $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  that we wish to implement using a given neural network. Our task boils down to selecting the interconnection weights in the network appropriately. As mentioned earlier, we refer to this task as *training* of the neural network or *learning* by the neural network. We use input-output examples of the given map to train the neural network. Specifically, let  $(\mathbf{x}_{d,1}, \mathbf{y}_{d,1}), \dots, (\mathbf{x}_{d,p}, \mathbf{y}_{d,p}) \in \mathbb{R}^n \times \mathbb{R}^m$ , where each  $\mathbf{y}_{d,i}$  is the output of the map  $\mathbf{F}$  corresponding to the input  $\mathbf{x}_{d,i}$ ; that is,  $\mathbf{y}_{d,i} = \mathbf{F}(\mathbf{x}_{d,i})$ . We refer to the set  $\{(\mathbf{x}_{d,1}, \mathbf{y}_{d,1}), \dots, (\mathbf{x}_{d,p}, \mathbf{y}_{d,p})\}$  as the *training set*. We train the neural network by adjusting the weights such that the map that is implemented by the network is close to the desired map  $\mathbf{F}$ . For this reason, we can think of neural networks as function approximators.

The form of learning described above can be thought of as learning with a teacher. The teacher supplies questions to the network in the form of  $\mathbf{x}_{d,1}, \dots, \mathbf{x}_{d,p}$  and tells the network the correct answers  $\mathbf{y}_{d,1}, \dots, \mathbf{y}_{d,p}$ . Training of the network then comprises applying a training algorithm that adjusts weights based on the error between the computed and desired outputs; that is, the difference between  $\mathbf{y}_{d,i} = \mathbf{F}(\mathbf{x}_{d,i})$  and the output of the neural network corresponding to  $\mathbf{x}_{d,i}$ . Having trained the neural network, our hope is that the network correctly generalizes the examples used in the training set. By this we mean that the network should correctly implement the mapping  $\mathbf{F}$  and produce the correct output corresponding to any input, including those that were not a part of the training set.

As we shall see in the remainder of this chapter, the training problem can be formulated as an optimization problem. We can then use optimization techniques and search methods (e.g., steepest descent, conjugate gradients [69], and quasi-Newton) for selection of the weights. The training algorithms are based on such optimization algorithms.

In the literature, for obvious reasons, the form of learning described above is referred to as *supervised learning*, a term which suggests that there is also a form of learning called *unsupervised learning*. Indeed, this is the case. However, unsupervised learning does not fit into the framework described above. Therefore, we do not discuss the idea of unsupervised learning any further. We refer the interested reader to [60].

## 13.2 Single-Neuron Training

Consider a single neuron, as shown in Figure 13.4. For this particular neuron, the activation function is simply the identity (linear function with unit slope). The neuron implements the following (linear) map from  $\mathbb{R}^n$  to  $\mathbb{R}$ :

$$y = \sum_{i=1}^n w_i x_i = \mathbf{x}^\top \mathbf{w},$$

where  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$  is the vector of inputs,  $y \in \mathbb{R}$  is the output, and  $\mathbf{w} = [w_1, \dots, w_n]^\top \in \mathbb{R}^n$  is the vector of weights. Suppose that we are given a map  $F : \mathbb{R}^n \rightarrow \mathbb{R}$ . We wish to find the value of the weights  $w_1, \dots, w_n$  such that the neuron approximates the map  $F$  as closely as possible. To do this, we use a training set consisting of  $p$  pairs  $\{(\mathbf{x}_{d,1}, y_{d,1}), \dots, (\mathbf{x}_{d,p}, y_{d,p})\}$ , where  $\mathbf{x}_{d,i} \in \mathbb{R}^n$  and  $y_{d,i} \in \mathbb{R}$ ,  $i = 1, \dots, p$ . For each  $i$ ,  $y_{d,i} = F(\mathbf{x}_{d,i})$  is the “desired” output corresponding to the given input  $\mathbf{x}_{d,i}$ . The training problem can then be formulated as the following optimization problem:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^p (y_{d,i} - \mathbf{x}_{d,i}^\top \mathbf{w})^2,$$

where the minimization is taken over all  $\mathbf{w} = [w_1, \dots, w_n]^\top \in \mathbb{R}^n$ . Note that the objective function represents the sum of the squared errors between the desired outputs  $y_{d,i}$  and the corresponding outputs of the neuron  $\mathbf{x}_{d,i}^\top \mathbf{w}$ . The factor of  $1/2$  is added for notational convenience and does not change the minimizer.

The objective function above can be written in matrix form as follows. First define the matrix  $\mathbf{X}_d \in \mathbb{R}^{n \times p}$  and vector  $\mathbf{y}_d \in \mathbb{R}^p$  by

$$\mathbf{X}_d = [\mathbf{x}_{d,1} \cdots \mathbf{x}_{d,p}],$$

$$\mathbf{y}_d = \begin{bmatrix} y_{d,1} \\ \vdots \\ y_{d,p} \end{bmatrix}.$$

Then, the optimization problem becomes

$$\text{minimize } \frac{1}{2} \|\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}\|^2.$$

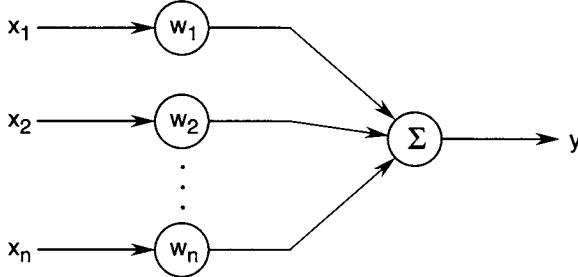


Figure 13.4 Single linear neuron.

There are two cases to consider in this optimization problem:  $p \leq n$  and  $p > n$ . We first consider the case where  $p \leq n$ , that is, where we have at most as many training pairs as the number of weights. For convenience, we assume that  $\text{rank } \mathbf{X}_d^\top = p$ . In this case there are an infinitely many points satisfying  $\mathbf{y}_d = \mathbf{X}_d^\top \mathbf{w}$ . Hence, there are infinitely many solutions to the optimization problem above, with the optimal objective function value of 0. Therefore, we have a choice of which optimal solution to select. A possible criterion for this selection is that of minimizing the solution norm. This is exactly the problem considered in Section 12.3. Recall that the minimum-norm solution is  $\mathbf{w}^* = \mathbf{X}_d(\mathbf{X}_d^\top \mathbf{X}_d)^{-1} \mathbf{y}_d$ . An efficient iterative algorithm for finding this solution is Kaczmarz's algorithm (discussed in Section 12.4). Kaczmarz's algorithm in this setting takes the form

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mu \frac{e_k \mathbf{x}_{d,R(k)+1}}{\|\mathbf{x}_{d,R(k)+1}\|^2},$$

where  $\mathbf{w}^{(0)} = \mathbf{0}$  and

$$e_k = y_{d,R(k)+1} - \mathbf{x}_{d,R(k)+1}^\top \mathbf{w}^{(k)}.$$

Recall that  $R(k)$  is the unique integer in  $\{0, \dots, p-1\}$  satisfying  $k = lp + R(k)$  for some integer  $l$ ; that is,  $R(k)$  is the remainder that results if we divide  $k$  by  $p$  (see Section 12.4 for further details on the algorithm).

The algorithm above was applied to the training of linear neurons by Widrow and Hoff (see [132] for some historical remarks). The single neuron together with the training algorithm above is illustrated in Figure 13.5 and is often called *Adaline*, an acronym for *adaptive linear element*.

We now consider the case where  $p > n$ . Here, we have more training points than the number of weights. We assume that  $\text{rank } \mathbf{X}_d^\top = n$ . In this case the objective function  $\frac{1}{2} \|\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}\|^2$  is simply a strictly convex quadratic function of  $\mathbf{w}$ , because the matrix  $\mathbf{X}_d \mathbf{X}_d^\top$  is a positive definite matrix. To solve this optimization problem, we have at our disposal the whole slew of unconstrained optimization algorithms considered in earlier chapters. For example, we can use a gradient algorithm, which in this case takes the form

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \alpha_k \mathbf{X}_d e^{(k)},$$

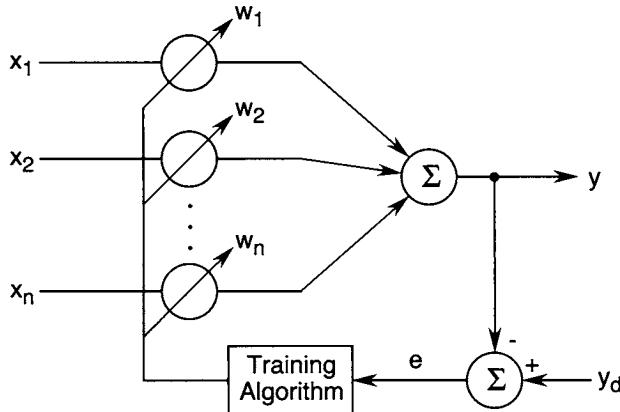


Figure 13.5 Adaline.

where  $e^{(k)} = y_d - \mathbf{X}_d^\top \mathbf{w}^{(k)}$ .

The discussion above assumed that the activation function for the neuron is the identity map. The derivation and analysis of the algorithms can be extended to the case of a general differentiable activation function  $f_a$ . Specifically, the output of the neuron in this case is given by

$$y = f_a \left( \sum_{i=1}^n w_i x_i \right) = f_a (\mathbf{x}^\top \mathbf{w}).$$

The algorithm for the case of a single training pair  $(\mathbf{x}_d, y_d)$  has the form

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mu \frac{e_k \mathbf{x}_d}{\|\mathbf{x}_d\|^2},$$

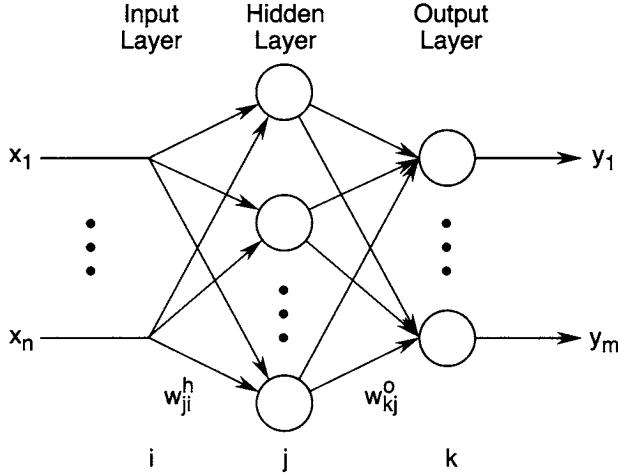
where the error is given by

$$e_k = y_d - f_a (\mathbf{x}_d^\top \mathbf{w}^{(k)}).$$

For a convergence analysis of the algorithm above, see [64].

### 13.3 The Backpropagation Algorithm

In Section 13.2 we considered the problem of training a single neuron. In this section we consider a neural network consisting of many layers. For simplicity of presentation, we restrict our attention to networks with three layers, as depicted in Figure 13.6. The three layers are referred to as the input, hidden, and output layers. There are  $n$  inputs  $x_i$ , where  $i = 1, \dots, n$ . We have  $m$  outputs  $y_s$ ,  $s = 1, \dots, m$ . There are  $l$  neurons in the hidden layer. The outputs



**Figure 13.6** Three-layered neural network.

of the neurons in the hidden layer are  $z_j$ , where  $j = 1, \dots, l$ . The inputs  $x_1, \dots, x_n$  are distributed to the neurons in the hidden layer. We may think of the neurons in the input layer as single-input-single-output linear elements, with each activation function being the identity map. In Figure 13.6 we do not explicitly depict the neurons in the input layer; instead, we illustrate the neurons as signal splitters. We denote the activation functions of the neurons in the hidden layer by  $f_j^h$ , where  $j = 1, \dots, l$ , and the activation functions of the neurons in the output layer by  $f_s^o$ , where  $s = 1, \dots, m$ . Note that each activation function is a function from  $\mathbb{R}$  to  $\mathbb{R}$ .

We denote the weights for inputs into the hidden layer by  $w_{ji}^h$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, l$ . We denote the weights for inputs from the hidden layer into the output layer by  $w_{sj}^o$ ,  $j = 1, \dots, l$ ,  $s = 1, \dots, m$ . Given the weights  $w_{ji}^h$  and  $w_{sj}^o$ , the neural network implements a map from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ . To find an explicit formula for this map, let us denote the input to the  $j$ th neuron in the hidden layer by  $v_j$  and the output of the  $j$ th neuron in the hidden layer by  $z_j$ . Then, we have

$$v_j = \sum_{i=1}^n w_{ji}^h x_i,$$

$$z_j = f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right).$$

The output from the  $s$ th neuron of the output layer is

$$y_s = f_s^o \left( \sum_{j=1}^l w_{sj}^o z_j \right).$$

Therefore, the relationship between the inputs  $x_i$ ,  $i = 1, \dots, n$ , and the  $s$ th output  $y_s$  is given by

$$\begin{aligned} y_s &= f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h(v_j) \right) \\ &= f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right) \right) \\ &= F_s(x_1, \dots, x_n). \end{aligned}$$

The overall mapping that the neural network implements is therefore given by

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} F_1(x_1, \dots, x_n) \\ \vdots \\ F_m(x_1, \dots, x_n) \end{bmatrix}.$$

We now consider the problem of training the neural network. As for the single neuron considered in Section 13.2, we analyze the case where the training set consists of a single pair  $(\mathbf{x}_d, \mathbf{y}_d)$ , where  $\mathbf{x}_d \in \mathbb{R}^n$  and  $\mathbf{y}_d \in \mathbb{R}^m$ . In practice, the training set consists of many such pairs, and training is typically performed with each pair at a time (see, e.g., [65] or [113]). Our analysis is therefore also relevant to the general training problem with multiple training pairs.

The training of the neural network involves adjusting the weights of the network such that the output generated by the network for the given input  $\mathbf{x}_d = [x_{d1}, \dots, x_{dn}]^\top$  is as close to  $\mathbf{y}_d$  as possible. Formally, this can be formulated as the following optimization problem:

$$\text{minimize } \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2,$$

where  $y_s$ ,  $s = 1, \dots, m$ , are the actual outputs of the neural network in response to the inputs  $x_{d1}, \dots, x_{dn}$ , as given by

$$y_s = f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right) \right).$$

This minimization is taken over all  $w_{ji}^h$ ,  $w_{sj}^o$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, l$ ,  $s = 1, \dots, m$ . For simplicity of notation, we use the symbol  $\mathbf{w}$  for the vector

$$\mathbf{w} = \{w_{ji}^h, w_{sj}^o : i = 1, \dots, n, j = 1, \dots, l, s = 1, \dots, m\}$$

and the symbol  $E$  for the objective function to be minimized; that is,

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2 \\ &= \frac{1}{2} \sum_{s=1}^m \left( y_{ds} - f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h \left( \sum_{i=1}^n w_{ji}^h x_{di} \right) \right) \right)^2. \end{aligned}$$

To solve the optimization problem above, we use a gradient algorithm with fixed step size. To formulate the algorithm, we need to compute the partial derivatives of  $E$  with respect to each component of  $\mathbf{w}$ . For this, let us first fix the indices  $i$ ,  $j$ , and  $s$ . We first compute the partial derivative of  $E$  with respect to  $w_{sj}^o$ . For this, we write

$$E(\mathbf{w}) = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - f_p^o \left( \sum_{q=1}^l w_{pq}^o z_q \right) \right)^2,$$

where for each  $q = 1, \dots, l$ ,

$$z_q = f_q^h \left( \sum_{i=1}^n w_{qi}^h x_{di} \right).$$

Using the chain rule, we obtain

$$\frac{\partial E}{\partial w_{sj}^o}(\mathbf{w}) = -(y_{ds} - y_s) f_s^{o'} \left( \sum_{q=1}^l w_{sq}^o z_q \right) z_j,$$

where  $f_s^{o'} : \mathbb{R} \rightarrow \mathbb{R}$  is the derivative of  $f_s^o$ . For simplicity of notation, we write

$$\delta_s = (y_{ds} - y_s) f_s^{o'} \left( \sum_{q=1}^l w_{sq}^o z_q \right).$$

We can think of each  $\delta_s$  as a scaled output error, because it is the difference between the actual output  $y_s$  of the neural network and the desired output  $y_{ds}$ , scaled by  $f_s^{o'} \left( \sum_{q=1}^l w_{sq}^o z_q \right)$ . Using the  $\delta_s$  notation, we have

$$\frac{\partial E}{\partial w_{sj}^o}(\mathbf{w}) = -\delta_s z_j.$$

We next compute the partial derivative of  $E$  with respect to  $w_{ji}^h$ . We start with the equation

$$E(\mathbf{w}) = \frac{1}{2} \sum_{p=1}^m \left( y_{dp} - f_p^o \left( \sum_{q=1}^l w_{pq}^o f_q^h \left( \sum_{r=1}^n w_{qr}^h x_{dr} \right) \right) \right)^2.$$

Using the chain rule once again, we get

$$\frac{\partial E}{\partial w_{ji}^h}(\mathbf{w}) = - \sum_{p=1}^m (y_{dp} - y_p) f_p^{o'} \left( \sum_{q=1}^l w_{pq}^o z_q \right) w_{pj}^o f_j^{h'} \left( \sum_{r=1}^n w_{jr}^h x_{dr} \right) x_{di},$$

where  $f_j^{h'} : \mathbb{R} \rightarrow \mathbb{R}$  is the derivative of  $f_j^h$ . Simplifying the above yields

$$\frac{\partial E}{\partial w_{ji}^h}(\mathbf{w}) = - \left( \sum_{p=1}^m \delta_p w_{pj}^o \right) f_j^{h'}(v_j) x_{di}.$$

We are now ready to formulate the gradient algorithm for updating the weights of the neural network. We write the update equations for the two sets of weights  $w_{sj}^o$  and  $w_{ji}^h$  separately. We have

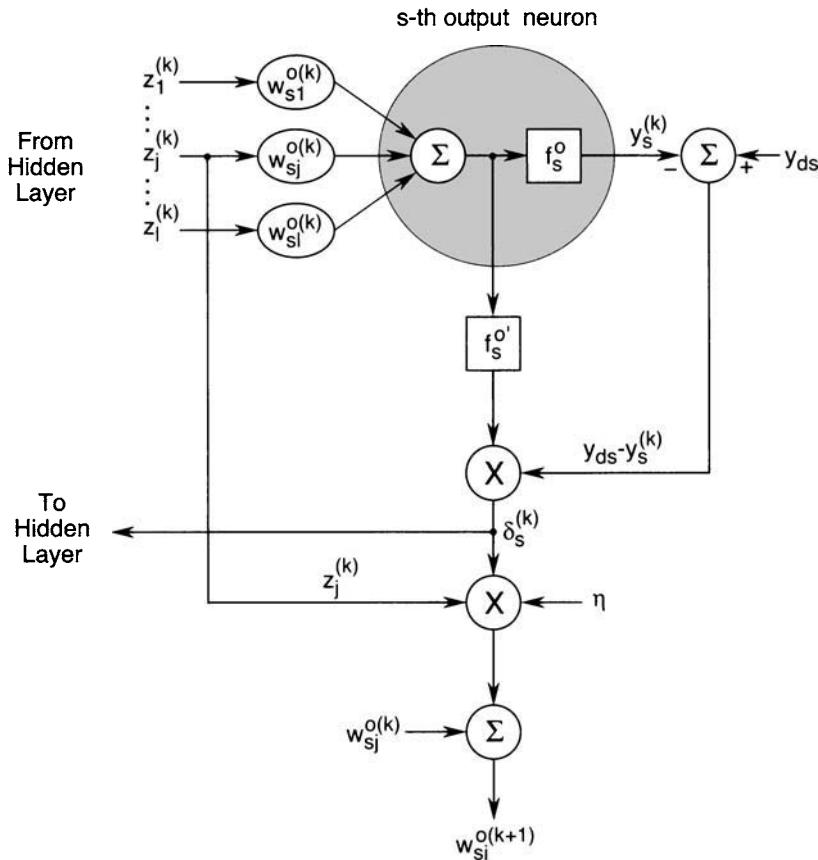
$$\begin{aligned} w_{sj}^{o(k+1)} &= w_{sj}^{o(k)} + \eta \delta_s^{(k)} z_j^{(k)}, \\ w_{ji}^{h(k+1)} &= w_{ji}^{h(k)} + \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^{o(k)} \right) f_j^{h'}(v_j^{(k)}) x_{di}, \end{aligned}$$

where  $\eta$  is the (fixed) step size and

$$\begin{aligned} v_j^{(k)} &= \sum_{i=1}^n w_{ji}^{h(k)} x_{di}, \\ z_j^{(k)} &= f_j^h(v_j^{(k)}), \\ y_s^{(k)} &= f_s^o \left( \sum_{q=1}^l w_{sq}^{o(k)} z_q^{(k)} \right), \\ \delta_s^{(k)} &= (y_{ds} - y_s^{(k)}) f_s^{o'} \left( \sum_{q=1}^l w_{sq}^{o(k)} z_q^{(k)} \right). \end{aligned}$$

The update equation for the weights  $w_{sj}^o$  of the output layer neurons is illustrated in Figure 13.7, whereas the update equation for the weights  $w_{ji}^h$  of the hidden layer neurons is illustrated in Figure 13.8.

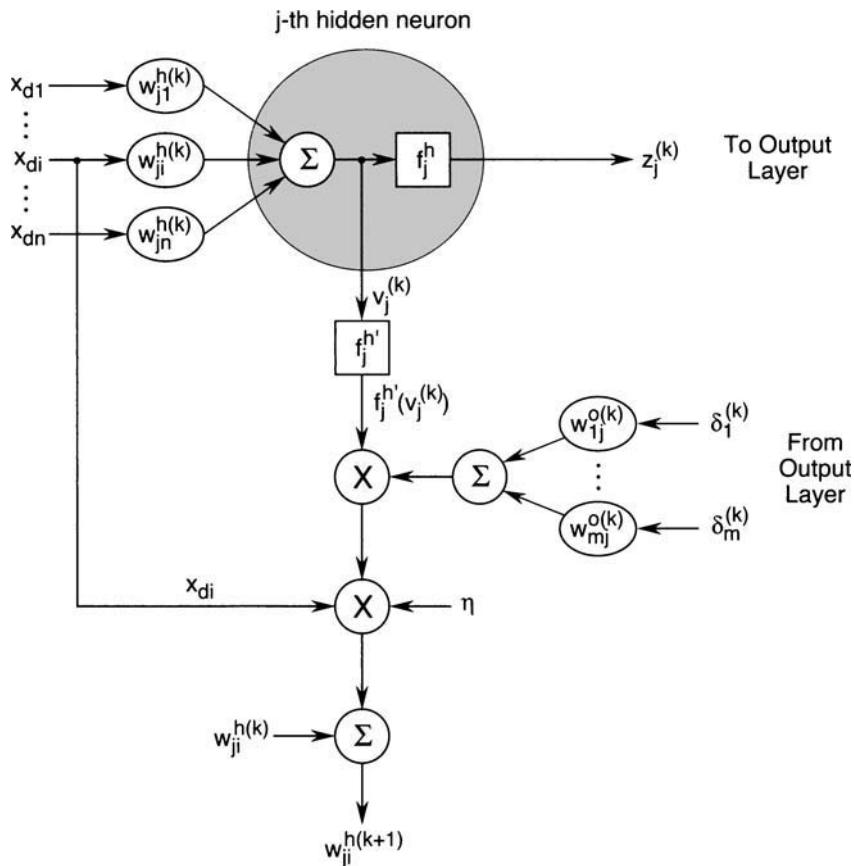
The update equations above are referred to in the literature as the *backpropagation algorithm*. The reason for the name *backpropagation* is that the output errors  $\delta_1^{(k)}, \dots, \delta_m^{(k)}$  are propagated back from the output layer to the hidden layer and are used in the update equation for the hidden layer weights, as illustrated in Figure 13.8. In the discussion above we assumed only a single hidden layer. In general, we may have multiple hidden layers—in this case the update equations for the weights will resemble the equations derived above. In the general case the output errors are propagated backward from layer to layer and are used to update the weights at each layer.



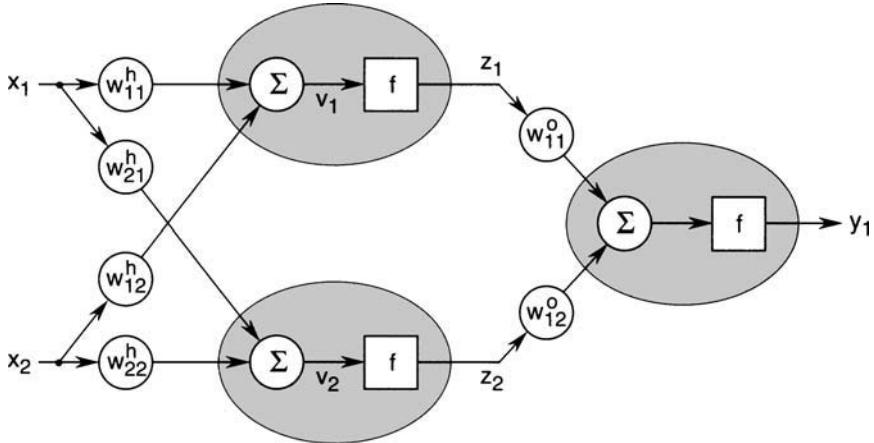
**Figure 13.7** Illustration of the update equation for the output layer weights.

We summarize the backpropagation algorithm qualitatively as follows. Using the inputs  $x_{di}$  and the current set of weights, we first compute the quantities  $v_j^{(k)}$ ,  $z_j^{(k)}$ ,  $y_s^{(k)}$ , and  $\delta_s^{(k)}$ , in turn. This is called the *forward pass* of the algorithm, because it involves propagating the input forward from the input layer to the output layer. Next, we compute the updated weights using the quantities computed in the forward pass. This is called the *reverse pass* of the algorithm, because it involves propagating the computed output errors  $\delta_s^{(k)}$  backward through the network. We illustrate the backpropagation procedure numerically in the following example.

**Example 13.1** Consider the simple feedforward neural network shown in Figure 13.9. The activation functions for all the neurons are given by  $f(v) = 1/(1 + e^{-v})$ . This particular activation function has the convenient property



**Figure 13.8** Illustration of the update equation for the hidden layer weights.



**Figure 13.9** Neural network for Example 13.1.

that  $f'(v) = f(v)(1 - f(v))$ . Therefore, using this property, we can write

$$\begin{aligned}\delta_1 &= (y_d - y_1)f' \left( \sum_{q=1}^2 w_{1q}^o z_q \right) \\ &= (y_d - y_1)f \left( \sum_{q=1}^2 w_{1q}^o z_q \right) \left( 1 - f \left( \sum_{q=1}^2 w_{1q}^o z_q \right) \right) \\ &= (y_d - y_1)y_1(1 - y_1).\end{aligned}$$

Suppose that the initial weights are  $w_{11}^{h(0)} = 0.1$ ,  $w_{12}^{h(0)} = 0.3$ ,  $w_{21}^{h(0)} = 0.3$ ,  $w_{22}^{h(0)} = 0.4$ ,  $w_{11}^{o(0)} = 0.4$ , and  $w_{12}^{o(0)} = 0.6$ . Let  $\mathbf{x}_d = [0.2, 0.6]^\top$  and  $y_d = 0.7$ . Perform one iteration of the backpropagation algorithm to update the weights of the network. Use a step size of  $\eta = 10$ .

To proceed, we first compute

$$\begin{aligned}v_1^{(0)} &= w_{11}^{h(0)}x_{d1} + w_{12}^{h(0)}x_{d2} = 0.2, \\ v_2^{(0)} &= w_{21}^{h(0)}x_{d1} + w_{22}^{h(0)}x_{d2} = 0.3.\end{aligned}$$

Next, we compute

$$\begin{aligned}z_1^{(0)} &= f(v_1^{(0)}) = \frac{1}{1 + e^{-0.2}} = 0.5498, \\ z_2^{(0)} &= f(v_2^{(0)}) = \frac{1}{1 + e^{-0.3}} = 0.5744.\end{aligned}$$

We then compute

$$y_1^{(0)} = f \left( w_{11}^{o(0)}z_1^{(0)} + w_{12}^{o(0)}z_2^{(0)} \right) = f(0.5646) = 0.6375,$$

which gives an output error of

$$\delta_1^{(0)} = (y_d - y_1^{(0)})y_1^{(0)}(1 - y_1^{(0)}) = 0.01444.$$

This completes the forward pass.

To update the weights, we use

$$\begin{aligned} w_{11}^{o(1)} &= w_{11}^{o(0)} + \eta\delta_1^{(0)}z_1^{(0)} = 0.4794, \\ w_{12}^{o(1)} &= w_{12}^{o(0)} + \eta\delta_1^{(0)}z_2^{(0)} = 0.6830, \end{aligned}$$

and, using the fact that  $f'(v_j^{(0)}) = f(v_j^{(0)})(1 - f(v_j^{(0)})) = z_j^{(0)}(1 - z_j^{(0)})$ , we get

$$\begin{aligned} w_{11}^{h(1)} &= w_{11}^{h(0)} + \eta\delta_1^{(0)}w_{11}^{o(0)}z_1^{(0)}(1 - z_1^{(0)})x_{d1} = 0.1029, \\ w_{12}^{h(1)} &= w_{12}^{h(0)} + \eta\delta_1^{(0)}w_{11}^{o(0)}z_1^{(0)}(1 - z_1^{(0)})x_{d2} = 0.3086, \\ w_{21}^{h(1)} &= w_{21}^{h(0)} + \eta\delta_1^{(0)}w_{12}^{o(0)}z_2^{(0)}(1 - z_2^{(0)})x_{d1} = 0.3042, \\ w_{22}^{h(1)} &= w_{22}^{h(0)} + \eta\delta_1^{(0)}w_{12}^{o(0)}z_2^{(0)}(1 - z_2^{(0)})x_{d2} = 0.4127. \end{aligned}$$

Thus, we have completed one iteration of the backpropagation algorithm. We can easily check that  $y_1^{(1)} = 0.6588$ , and hence  $|y_d - y_1^{(1)}| < |y_d - y_1^{(0)}|$ ; that is, the actual output of the neural network has become closer to the desired output as a result of updating the weights.

After 15 iterations of the backpropagation algorithm, we get

$$\begin{aligned} w_{11}^{o(15)} &= 0.6365, \\ w_{12}^{o(15)} &= 0.8474, \\ w_{11}^{h(15)} &= 0.1105, \\ w_{12}^{h(15)} &= 0.3315, \\ w_{21}^{h(15)} &= 0.3146, \\ w_{22}^{h(15)} &= 0.4439. \end{aligned}$$

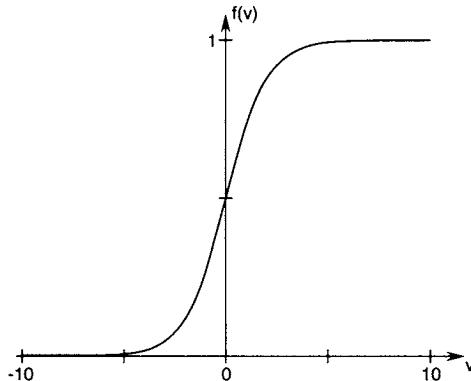
The resulting value of the output corresponding to the input  $\mathbf{x}_d = [0.2, 0.6]^\top$  is  $y_1^{(15)} = 0.6997$ . ■

In the example above, we considered an activation function of the form

$$f(v) = \frac{1}{1 + e^{-v}}.$$

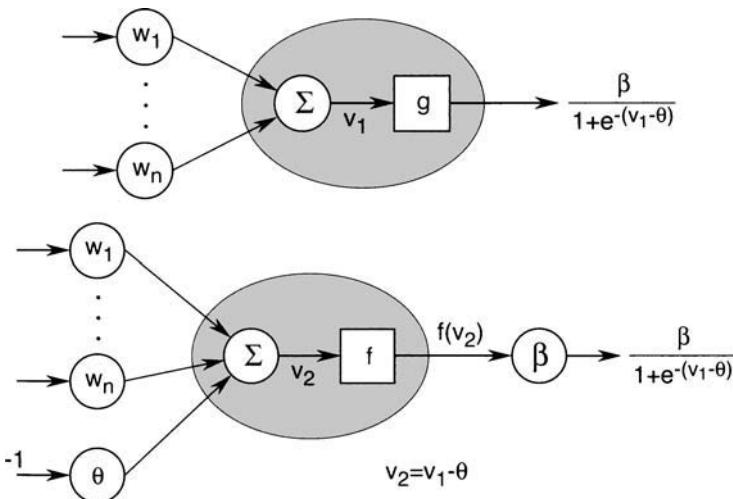
This function is called a *sigmoid* and is a popular activation function used in practice. The sigmoid function is illustrated in Figure 13.10. It is possible to use a more general version of the sigmoid function, of the form

$$g(v) = \frac{\beta}{1 + e^{-(v-\theta)}}.$$



**Figure 13.10** Sigmoid function.

The parameters  $\beta$  and  $\theta$  represent *scale* and *shift* (or *location*) *parameters* respectively. The parameter  $\theta$  is often interpreted as a threshold. If such an activation function is used in a neural network, we would also want to adjust the values of the parameters  $\beta$  and  $\theta$ , which also affect the value of the objective function to be minimized. However, it turns out that these parameters can be incorporated into the backpropagation algorithm simply by treating them as additional weights in the network. Specifically, we can represent a neuron with activation function  $g$  as one with activation function  $f$  with the addition of two extra weights, as shown in Figure 13.11.



**Figure 13.11** Two configurations that are equivalent.

**Example 13.2** Consider the same neural network as in Example 13.1. We introduce shift parameters  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$  to the activation functions in the neurons. Using the configuration illustrated in Figure 13.11, we can incorporate the shift parameters into the backpropagation algorithm. We have

$$\begin{aligned} v_1 &= w_{11}^h x_{d1} + w_{12}^h x_{d2} - \theta_1, \\ v_2 &= w_{21}^h x_{d1} + w_{22}^h x_{d2} - \theta_2, \\ z_1 &= f(v_1), \\ z_2 &= f(v_2), \\ y_1 &= f(w_{11}^o z_1 + w_{12}^o z_2 - \theta_3), \\ \delta_1 &= (y_d - y_1)y_1(1 - y_1), \end{aligned}$$

where  $f$  is the sigmoid function:

$$f(v) = \frac{1}{1 + e^{-v}}.$$

The components of the gradient of the objective function  $E$  with respect to the shift parameters are

$$\begin{aligned} \frac{\partial E}{\partial \theta_1}(\mathbf{w}) &= \delta_1 w_{11}^o z_1 (1 - z_1), \\ \frac{\partial E}{\partial \theta_2}(\mathbf{w}) &= \delta_1 w_{12}^o z_2 (1 - z_2), \\ \frac{\partial E}{\partial \theta_3}(\mathbf{w}) &= \delta_1. \end{aligned}$$

■

In the next example, we apply the network discussed in Example 13.2 to solve the celebrated *exclusive OR (XOR) problem* (see [113]).

**Example 13.3** Consider the neural network of Example 13.2. We wish to train the neural network to approximate the *exclusive OR (XOR)* function, defined in Table 13.1. Note that the XOR function has two inputs and one output.

To train the neural network, we use the following training pairs:

$$\begin{aligned} \mathbf{x}_{d,1} &= [0, 0]^\top, & y_{d,1} &= 0, \\ \mathbf{x}_{d,2} &= [0, 1]^\top, & y_{d,2} &= 1, \\ \mathbf{x}_{d,3} &= [1, 0]^\top, & y_{d,3} &= 1, \\ \mathbf{x}_{d,4} &= [1, 1]^\top, & y_{d,4} &= 0. \end{aligned}$$

We now apply the backpropagation algorithm to train the network using the training pairs above. To do this, we apply one pair per iteration in a cyclic

**Table 13.1** Truth Table for XOR Function

$x_1$	$x_2$	$F(x_1, x_2)$
0	0	0
0	1	1
1	0	1
1	1	0

**Table 13.2** Response of the Trained Network of Example 13.3

$x_1$	$x_2$	$y_1$
0	0	0.007
0	1	0.99
1	0	0.99
1	1	0.009

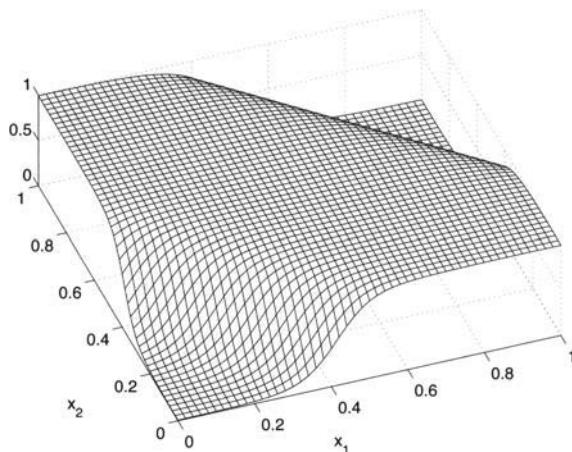
fashion. In other words, in the  $k$ th iteration of the algorithm, we apply the pair  $(\mathbf{x}_{d,R(k)+1}, y_{d,R(k)+1})$ , where, as in Kaczmarz's algorithm,  $R(k)$  is the unique integer in  $\{0, 1, 2, 3\}$  satisfying  $k = 4l + R(k)$  for some integer  $l$ ; that is,  $R(k)$  is the remainder that results if we divide  $k$  by 4 (see Section 12.4).

The experiment yields the following weights (see Exercise 13.5):

$$\begin{aligned} w_{11}^o &= -11.01, \\ w_{12}^o &= 10.92, \\ w_{11}^h &= -7.777, \\ w_{12}^h &= -8.403, \\ w_{21}^h &= -5.593, \\ w_{22}^h &= -5.638, \\ \theta_1 &= -3.277, \\ \theta_2 &= -8.357, \\ \theta_3 &= 5.261. \end{aligned}$$

Table 13.2 shows the output of the neural network with the weights above corresponding to the training input data. Figure 13.12 shows a plot of the function that is implemented by this neural network. ■

For a more comprehensive treatment of neural networks, see [58], [59], or [137]. For applications of neural networks to optimization, signal processing, and control problems, see [28] and [67].



**Figure 13.12** Plot of the function implemented by the trained network of Example 13.3.

## EXERCISES

**13.1** Consider a single linear neuron, with  $n$  inputs (see Figure 13.4). Suppose that we are given  $\mathbf{X}_d \in \mathbb{R}^{n \times p}$  and  $\mathbf{y}_d \in \mathbb{R}^p$  representing  $p$  training pairs, where  $p > n$ . The objective function to be minimized in the training of the neuron is

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}\|^2.$$

- a. Find the gradient of the objective function.
- b. Write the conjugate gradient algorithm for training the neuron.
- c. Suppose that we wish to approximate the function  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$F(\mathbf{x}) = (\sin x_1)(\cos x_2).$$

Use the conjugate gradient algorithm from part b to train the linear neuron, using the following training points:

$$\{\mathbf{x} : x_1, x_2 = -0.5, 0, 0.5\}.$$

It may helpful to use the MATLAB program from Exercise 10.11.

- d. Plot the level sets of the objective function for the problem in part c, at levels 0.01, 0.1, 0.2, and 0.4. Check if the solution in part c agrees with the level sets.

- e. Plot the error function  $e(\mathbf{x}) = F(\mathbf{x}) - \mathbf{w}^{*\top} \mathbf{x}$  versus  $x_1$  and  $x_2$ , where  $\mathbf{w}^*$  is the optimal weight vector obtained in part c.

**13.2** Consider the Adaline, depicted in Figure 13.5. Assume that we have a single training pair  $(\mathbf{x}_d, y_d)$ , where  $\mathbf{x}_d \neq \mathbf{0}$ . Suppose that we use the Widrow-Hoff algorithm to adjust the weights:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mu \frac{e_k \mathbf{x}_d}{\mathbf{x}_d^\top \mathbf{x}_d},$$

where  $e_k = y_d - \mathbf{x}_d^\top \mathbf{w}^{(k)}$ .

- a. Write an expression for  $e_{k+1}$  as a function of  $e_k$  and  $\mu$ .
- b. Find the largest range of values for  $\mu$  for which  $e_k \rightarrow 0$  (for any initial condition  $\mathbf{w}^{(0)}$ ).

**13.3** As in Exercise 13.2, consider the Adaline. Consider the case in which there are multiple pairs in the training set  $\{(\mathbf{x}_{d,1}, y_{d,1}), \dots, (\mathbf{x}_{d,p}, y_{d,p})\}$ , where  $p \leq n$  and  $\text{rank } \mathbf{X}_d = p$  (the matrix  $\mathbf{X}_d$  has  $\mathbf{x}_{d,i}$  as its  $i$ th column). Suppose that we use the following training algorithm:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{X}_d (\mathbf{X}_d^\top \mathbf{X}_d)^{-1} \boldsymbol{\mu} \mathbf{e}^{(k)},$$

where  $\mathbf{e}^{(k)} = \mathbf{y}_d - \mathbf{X}_d^\top \mathbf{w}^{(k)}$  and  $\boldsymbol{\mu}$  is a given constant  $p \times p$  matrix.

- a. Find an expression for  $\mathbf{e}^{(k+1)}$  as a function of  $\mathbf{e}^{(k)}$  and  $\boldsymbol{\mu}$ .
- b. Find a necessary and sufficient condition on  $\boldsymbol{\mu}$  for which  $\mathbf{e}^{(k)} \rightarrow \mathbf{0}$  (for any initial condition  $\mathbf{w}^{(0)}$ ).

**13.4** Consider the three-layered neural network described in Example 13.1 (see Figure 13.9). Implement the backpropagation algorithm for this network in MATLAB. Test the algorithm for the training pair  $\mathbf{x}_d = [0, 1]^\top$  and  $y_d = 1$ . Use a step size of  $\eta = 50$  and initial weights as in the Example 13.1.

**13.5** Consider the neural network of Example 13.3, with training pairs for the XOR problem. Use MATLAB to implement the training algorithm described in Example 13.3, with a step size of  $\eta = 10.0$ . Tabulate the outputs of the trained network corresponding to the training input data.



# CHAPTER 14

---

## GLOBAL SEARCH ALGORITHMS

---

### 14.1 Introduction

The iterative algorithms in previous chapters, in particular gradient methods, Newton's method, conjugate gradient methods, and quasi-Newton methods, start with an initial point and then generate a sequence of iterates. Typically, the best we can hope for is that the sequence converges to a local minimizer. For this reason, it is often desirable for the initial point to be close to a global minimizer. Moreover, these methods require first derivatives (and also second derivatives in the case of Newton's method).

In this chapter we discuss various search methods that are global in nature in the sense that they attempt to search throughout the entire feasible set. These methods use only objective function values and do not require derivatives. Consequently, they are applicable to a much wider class of optimization problems. In some cases, they can also be used to generate “good” initial (starting) points for the iterative methods discussed in earlier chapters. Some of the methods we discuss in this chapter (specifically, the randomized search methods) are also used in combinatorial optimization, where the feasible set is finite (discrete), but typically large.

## 14.2 The Nelder-Mead Simplex Algorithm

The method originally proposed by Spendley, Hext, and Hinsworth [122] in 1962 was improved by Nelder and Mead [97] in 1965 and it is now commonly referred to as the *Nelder-Mead simplex algorithm*. A contemporary view of the algorithm is provided in the well-written paper by Lagarias et al. [82]. In our exposition, we use the notation of this paper.

The Nelder-Mead algorithm is a derivative-free method. The method uses the concept of a simplex. A *simplex* is a geometric object determined by an assembly of  $n + 1$  points,  $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_n$ , in the  $n$ -dimensional space such that

$$\det \begin{bmatrix} \mathbf{p}_0 & \mathbf{p}_1 & \cdots & \mathbf{p}_n \\ 1 & 1 & \cdots & 1 \end{bmatrix} \neq 0.$$

This condition ensures that two points in  $\mathbb{R}$  do not coincide, three points in  $\mathbb{R}^2$  are not collinear, four points in  $\mathbb{R}^3$  are not coplanar, and so on. Thus, simplex in  $\mathbb{R}$  is a line segment, in  $\mathbb{R}^2$  it is a triangle, while a simplex in  $\mathbb{R}^3$  is a tetrahedron; in each case it encloses a finite  $n$ -dimensional volume.

Suppose that we wish to minimize  $f(\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^n$ . To start the algorithm, we initialize a simplex of  $n + 1$  points. A possible way to set up a simplex, as suggested by Jang, Sun, and Mizutani [67], is to start with an initial point  $\mathbf{x}^{(0)} = \mathbf{p}_0$  and generate the remaining points of the initial simplex as follows:

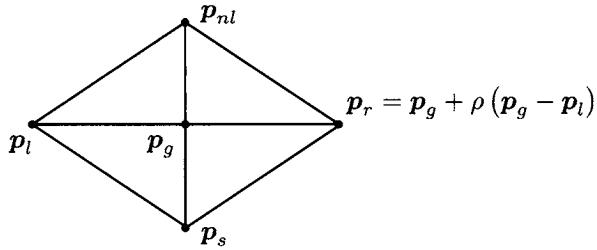
$$\mathbf{p}_i = \mathbf{p}_0 + \lambda_i \mathbf{e}_i, \quad i = 1, 2, \dots, n,$$

where the  $\mathbf{e}_i$  are unit vectors constituting the natural basis of  $\mathbb{R}^n$  as described in Section 2.1. The positive constant coefficients  $\lambda_i$  are selected in such a way that their magnitudes reflect the length scale of the optimization problem. Our objective is to modify the initial simplex stage by stage so that the resulting simplices converge toward the minimizer. At each iteration we evaluate the function  $f$  at each point of the simplex. In the function minimization process, the point with the largest function value is replaced with another point. The process for modifying the simplex continues until it converges toward the function minimizer.

We now present the rules for modifying the simplex stage by stage. To aid in our presentation, we use a two-dimensional example to illustrate the rules. We begin by selecting the initial set of  $n + 1$  points that are to form the initial simplex. We next evaluate  $f$  at each point and order the  $n + 1$  vertices to satisfy

$$f(\mathbf{p}_0) \leq f(\mathbf{p}_1) \leq \cdots \leq f(\mathbf{p}_n).$$

For the two-dimensional case we let  $\mathbf{p}_l$ ,  $\mathbf{p}_{nl}$ , and  $\mathbf{p}_s$  denote the points of the simplex for which  $f$  is largest, next largest, and smallest; that is, because we wish to minimize  $f$ , the vertex  $\mathbf{p}_s$  is the best vertex,  $\mathbf{p}_l$  is the worst vertex, and  $\mathbf{p}_{nl}$  is the next-worst vertex. We next compute  $\mathbf{p}_g$ , the centroid (center



**Figure 14.1** Reflecting  $\mathbf{p}_l$  in  $\mathbf{p}_g$  with a reflection coefficient  $\rho$ .

of gravity) of the best  $n$  points:

$$\mathbf{p}_g = \sum_{i=0}^{n-1} \frac{\mathbf{p}_i}{n}.$$

In our two-dimensional case,  $n = 2$ , we would have

$$\mathbf{p}_g = \frac{1}{2} (\mathbf{p}_{nl} + \mathbf{p}_s).$$

We then reflect the worst vertex,  $\mathbf{p}_l$ , in  $\mathbf{p}_g$  using a reflection coefficient  $\rho > 0$  to obtain the reflection point

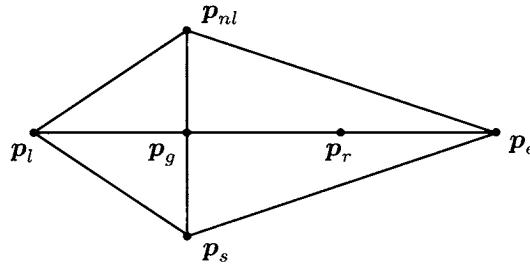
$$\mathbf{p}_r = \mathbf{p}_g + \rho (\mathbf{p}_g - \mathbf{p}_l).$$

A typical value is  $\rho = 1$ . The operation above is illustrated in Figure 14.1. We proceed to evaluate  $f$  at  $\mathbf{p}_r$  to obtain  $f_r = f(\mathbf{p}_r)$ . If  $f_0 \leq f_r < f_{n-1}$  [i.e., if  $f_r$  lies between  $f_s = f(\mathbf{p}_s)$  and  $f_{nl} = f(\mathbf{p}_{nl})$ ], then the point  $\mathbf{p}_r$  replaces  $\mathbf{p}_l$  to form a new simplex, and we terminate the iteration. We proceed to repeat the process. Thus, we compute the centroid of the best  $n$  vertices of the new simplex and again reflect the point with the largest function  $f$  value in the centroid obtained for the best  $n$  points of the new simplex.

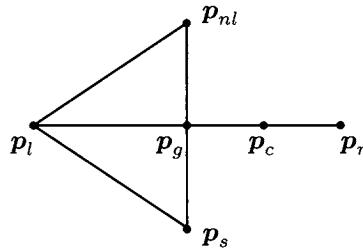
If, however,  $f_r < f_s = f_0$ , so that the point  $\mathbf{p}_r$  yields the smallest function value among the points of the simplex, we argue that this direction is a good one. In this case we increase the distance traveled using an *expansion coefficient*  $\chi > 1$  (e.g.,  $\chi = 2$ ) to obtain

$$\mathbf{p}_e = \mathbf{p}_g + \chi (\mathbf{p}_r - \mathbf{p}_g).$$

The operation above yields a new point on the line  $\mathbf{p}_l\mathbf{p}_g\mathbf{p}_r$  extended beyond  $\mathbf{p}_r$ . We illustrate this operation in Figure 14.2. If  $f_e < f_r$  now, the expansion is declared a success and  $\mathbf{p}_e$  replaces  $\mathbf{p}_l$  in the next simplex. If, on the other hand,  $f_e \geq f_r$ , the expansion is a failure and  $\mathbf{p}_r$  replaces  $\mathbf{p}_l$ .



**Figure 14.2** Expansion operation with the expansion coefficient  $\chi$ .



**Figure 14.3** Outside contraction operation for the case when  $f_r \in [f_{nl}, f_l)$ .

Finally, if  $f_r \geq f_{nl}$ , the reflected point  $p_r$  would constitute the point with the largest function value in the new simplex. Then in the next step it would be reflected in  $p_g$ , probably an unfruitful operation. Instead, this case is dealt with by a *contraction* operation in one of two ways. First, if  $f_r \geq f_{nl}$  and  $f_r < f_l$ , then we contract  $(p_r - p_g)$  with a contraction coefficient  $0 < \gamma < 1$  (e.g.,  $\gamma = 1/2$ ) to obtain

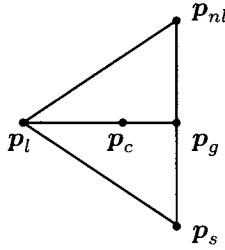
$$\mathbf{p}_c = \mathbf{p}_g + \gamma (\mathbf{p}_r - \mathbf{p}_g).$$

We refer to this operation as the *outside contraction*. See Figure 14.3 for an illustration of this operation. If, on the other hand,  $f_r \geq f_{nl}$  and  $f_r \geq f_l$ , then  $p_l$  replaces  $p_r$  in the contraction operation and we get

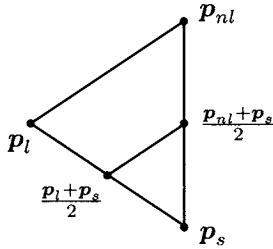
$$\mathbf{p}_c = \mathbf{p}_g + \gamma (\mathbf{p}_l - \mathbf{p}_g).$$

This operation, referred to as the *inside contraction*, is illustrated in Figure 14.4.

If, in either case,  $f_c \leq f_l$ , the contraction is considered a success, and we replace  $p_l$  with  $p_c$  in the new simplex. If, however,  $f_c > f_l$ , the contraction



**Figure 14.4** Inside contraction operation for the case when  $f_r \geq f_l$ .



**Figure 14.5** Shrinkage operation.

is a failure, and in this case a new simplex can be formed by retaining  $p_s$  only and halving the distance from  $p_s$  to every other point in the simplex. We can refer to this event as a shrinkage operation. The shrinkage operation is illustrated in Figure 14.5. In general, the shrink step produces the  $n$  new vertices of the new simplex according to the formula

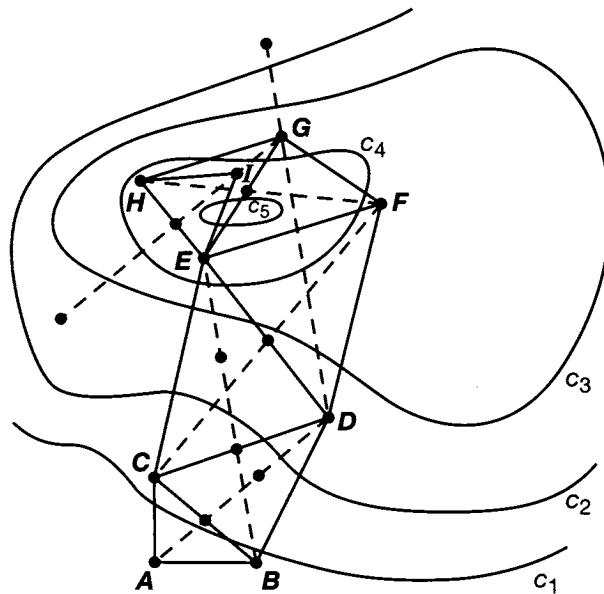
$$\mathbf{v}_i = \mathbf{p}_s + \sigma(\mathbf{p}_i - \mathbf{p}_s), \quad i = 1, 2, \dots, n,$$

where  $\sigma = 1/2$ . Hence, the vertices of the new simplex are  $\mathbf{p}_s, \mathbf{v}_1, \dots, \mathbf{v}_n$ .

When implementing the simplex algorithm, we need a tie-breaking rule to order points in the case of equal function values. Lagarias et al. [82] propose tie-breaking rules that assign to the new vertex the highest possible index consistent with the relation

$$f(\mathbf{p}_0) \leq f(\mathbf{p}_1) \leq \cdots \leq f(\mathbf{p}_n).$$

In Figure 14.6 we illustrate the simplex search method by showing the first few stages of the search for a minimizer of a function of two variables. Our



**Figure 14.6** The simplex search method applied to minimization of a function of two variables.

drawing is inspired by a figure in Layton [84, p. 225]. The starting simplex is composed of the vertices  $A$ ,  $B$ , and  $C$ . The vertices  $D$  and  $E$  are obtained by the expansion operation. The vertex  $F$  is obtained by the reflection operation. The vertex  $G$  is obtained using the outside contraction operation, while the vertex  $I$  is obtained employing the inside contraction operation. For the sake of clarity we terminate the process with the simplex composed of the vertices  $E$ ,  $H$ , and  $I$ . The process may, of course, be continued beyond this simplex.

We add that a variant of the simplex method described above is presented in Jang et al. [67], where they use the centroid of the entire simplex rather than the centroid of the best  $n$  vertices of the simplex. That is, Jang et al. [67] compute the point  $p_g$  using the  $n+1$  vertices of the simplex. In addition, they use only the inside contraction and they do not use the outside contraction operation.

### 14.3 Simulated Annealing

#### Randomized Search

Simulated annealing is an instance of a randomized search method. A *randomized search method*, also sometimes called a *probabilistic search method*, is an algorithm that searches the feasible set of an optimization problem by

considering randomized samples of candidate points in the set. The simulated annealing algorithm was first suggested for optimization by Kirkpatrick et al. [75] based on techniques of Metropolis et al. [91]. An early application to image processing was described by Geman and Geman [48].

As usual, suppose that we wish to solve an optimization problem of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

The basic assumption in randomized search is our ability to select a random sample from the feasible set  $\Omega$ . Typically, we start a randomized search process by selecting a random initial point  $\mathbf{x}^{(0)} \in \Omega$ . Then, we select a random next-candidate point, usually close to  $\mathbf{x}^{(0)}$ .

More formally, we assume that for any  $\mathbf{x} \in \Omega$ , there is a set  $N(\mathbf{x}) \subset \Omega$  such that we can generate a random sample from this set. Typically,  $N(\mathbf{x})$  is a set of points that are “close” to  $\mathbf{x}$ , and for this reason we usually think of  $N(\mathbf{x})$  as a “neighborhood” of  $\mathbf{x}$  [we use the term *neighborhood* for  $N(\mathbf{x})$  even in the general case where the points in it are arbitrary, not necessarily close to  $\mathbf{x}$ ]. When we speak of generating a random point in  $N(\mathbf{x})$ , we mean that there is a prespecified distribution over  $N(\mathbf{x})$ , and we sample a point with this distribution. Often, this distribution is chosen to be uniform over  $N(\mathbf{x})$ ; other distributions are also used, including Gaussian and Cauchy.

Before discussing the simulated annealing method, we first consider a simple randomized search algorithm, which we will call *naive random search*.

### Naive Random Search Algorithm

1. Set  $k := 0$ . Select an initial point  $\mathbf{x}^{(0)} \in \Omega$ .
2. Pick a candidate point  $\mathbf{z}^{(k)}$  at random from  $N(\mathbf{x}^{(k)})$ .
3. If  $f(\mathbf{z}^{(k)}) < f(\mathbf{x}^{(k)})$ , then set  $\mathbf{x}^{(k+1)} = \mathbf{z}^{(k)}$ ; else, set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ .
4. If stopping criterion satisfied, then stop.
5. Set  $k := k + 1$ , go to step 2.

Note that the algorithm above has the familiar form  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$ , where  $\mathbf{d}^{(k)}$  is randomly generated. By design, the direction  $\mathbf{d}^{(k)}$  either is  $\mathbf{0}$  or is a descent direction. Typical stopping criteria include reaching a certain number of iterations, or reaching a certain objective function value.

### Simulated Annealing Algorithm

The main problem with the naive random search method is that it may get stuck in a region around a local minimizer. This is easy to imagine; for

example, if  $\mathbf{x}^{(0)}$  is a local minimizer and  $N(\mathbf{x}^{(0)})$  is sufficiently small that all points in it have no smaller objective function value than  $\mathbf{x}^{(0)}$  [i.e.,  $\mathbf{x}^{(0)}$  is a global minimizer of  $f$  over  $N(\mathbf{x}^{(0)})$ ], then clearly the algorithm will be stuck and will never find a point outside of  $N(\mathbf{x}^{(0)})$ . To prevent getting stuck in a region around a local minimizer, we need a way to consider points outside this region. One way to achieve this goal is to make sure that at each  $k$ , the neighborhood  $N(\mathbf{x}^{(k)})$  is a very large set. Indeed, if  $N(\mathbf{x}^{(k)})$  is sufficiently large, then we are guaranteed that the algorithm will converge (in some sense) to a global minimizer. An extreme example of this case is where  $N(\mathbf{x}) = \Omega$  for any  $\mathbf{x} \in \Omega$  (in this case running  $k$  iterations of the naive random search algorithm amounts to finding the best point among  $k$  randomly chosen points in  $\Omega$ ). However, having too large a neighborhood in the search algorithm results in a slow search process, because the sampling of candidate points to consider is spread out, making it more unlikely to find a better candidate point.

Another way to overcome the problem of getting stuck in a region around a local minimizer is to modify the naive search algorithm so that we can “climb out” of such a region. This means that the algorithm may accept a new point that is *worse* than the current point. The simulated annealing algorithm incorporates such a mechanism.

### Simulated Annealing Algorithm

1. Set  $k := 0$ ; select an initial point  $\mathbf{x}^{(0)} \in \Omega$ .
2. Pick a candidate point  $\mathbf{z}^{(k)}$  at random from  $N(\mathbf{x}^{(k)})$ .
3. Toss a coin with probability of HEAD equal to  $p(k, f(\mathbf{z}^{(k)}), f(\mathbf{x}^{(k)}))$ . If HEAD, then set  $\mathbf{x}^{(k+1)} = \mathbf{z}^{(k)}$ ; else, set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ .
4. If the stopping criterion is satisfied, then stop.
5. Set  $k := k + 1$ , go to step 2.

In step 3, the use of a “coin toss” is simply descriptive for a randomized decision—we do not mean literally that an actual coin needs to be tossed.

As in naive random search, the simulated annealing algorithm above has the familiar form  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$ , where  $\mathbf{d}^{(k)}$  is randomly generated. But in simulated annealing the direction  $\mathbf{d}^{(k)}$  might be an ascent direction. However, as the algorithm progresses, we can keep track of the *best-so-far point*—this is a point  $\mathbf{x}_{\text{best}}^{(k)}$  which, at each  $k$ , is equal to a  $\mathbf{x}^{(j)}$ ,  $j \in \{0, \dots, k\}$ , such that  $f(\mathbf{x}^{(j)}) \leq f(\mathbf{x}^{(i)})$  for all  $i \in \{0, \dots, k\}$ . The best-so-far point can be updated at each step  $k$  as follows:

$$\mathbf{x}_{\text{best}}^{(k)} = \begin{cases} \mathbf{x}^{(k)} & \text{if } f(\mathbf{x}^{(k)}) < f(\mathbf{x}_{\text{best}}^{(k-1)}) \\ \mathbf{x}_{\text{best}}^{(k-1)} & \text{otherwise.} \end{cases}$$

By keeping track of the best-so-far point, we can treat the simulated annealing algorithm simply as a search procedure; the best-so-far point is what we eventually use when the algorithm stops. This comment applies not only to simulated annealing, but other search techniques as well (including the randomized algorithms presented in the next two sections).

The major difference between simulated annealing and naive random search is that in step 5, there is some probability that we set the next iterate to be equal to the random point selected from the neighborhood, even if that point turns out to be worse than the current iterate. This probability is called the *acceptance probability*. For the algorithm to work properly, the acceptance probability must be chosen appropriately. A typical choice is

$$p(k, f(\mathbf{z}^{(k)}), f(\mathbf{x}^{(k)})) = \min\{1, \exp(-(f(\mathbf{z}^{(k)}) - f(\mathbf{x}^{(k)}))/T_k)\},$$

where  $\exp$  is the exponential function and  $T_k$  represents a positive sequence called the *temperature schedule* or *cooling schedule*. This form of acceptance probability is usually credited to Boltzmann and leads to a simulated annealing algorithm that behaves as a Gibbs sampler (a method of probabilistic sampling based on the Gibbs distribution).

Notice that if  $f(\mathbf{z}^{(k)}) \leq f(\mathbf{x}^{(k)})$ , then  $p(k, f(\mathbf{z}^{(k)}), f(\mathbf{x}^{(k)})) = 1$ , which means that we set  $\mathbf{x}^{(k+1)} = \mathbf{z}^{(k)}$  (i.e., we move to the point  $\mathbf{z}^{(k)}$ ). However, if  $f(\mathbf{z}^{(k)}) > f(\mathbf{x}^{(k)})$ , there is still a positive probability of setting  $\mathbf{x}^{(k+1)} = \mathbf{z}^{(k)}$ ; this probability is equal to

$$\exp\left(-\frac{f(\mathbf{z}^{(k)}) - f(\mathbf{x}^{(k)})}{T_k}\right).$$

Note that the larger the difference between  $f(\mathbf{z}^{(k)})$  and  $f(\mathbf{x}^{(k)})$ , the less likely we are to move to the worse point  $\mathbf{z}^{(k)}$ . Similarly, the smaller the value of  $T_k$ , the less likely we are to move to  $\mathbf{z}^{(k)}$ . It is typical to let the “temperature”  $T_k$  be monotonically decreasing to 0 (hence the word *cooling*). In other words, as the iteration index  $k$  increases, the algorithm becomes increasingly reluctant to move to a worse point. The intuitive reason for this behavior is that initially we wish to actively explore the feasible set, but with time we would like to be less active in exploration so that we spend more time in a region around a global minimizer. In other words, the desired behavior is this: Initially, the algorithm jumps around and is more likely to climb out of regions around local minimizers, but with time it settles down and is more likely to spend time around a global minimizer.

The term *annealing* comes from the field of metallurgy, where it refers to a technique for improving the property of metals. The basic procedure is to heat up a piece of metal and then cool it down in a controlled fashion. When the metal is first heated, the atoms in it become unstuck from their initial positions (with some level of internal energy). Then, as cooling takes place, the atoms gradually configure themselves in states of lower internal energy. Provided that the cooling is sufficiently slow, the final internal energy is lower

than the initial internal energy, thereby refining the crystalline structure and reducing defects.

In an analogous way, the temperature in simulated annealing must be cooled in a controlled fashion. In particular, the cooling should be sufficiently slow. In a seminal paper, Hajek [56] provides a rigorous analysis of the cooling schedule for convergence of the algorithm to a global minimizer. Specifically, he shows that an appropriate cooling schedule is

$$T_k = \frac{\gamma}{\log(k+2)},$$

where  $\gamma > 0$  is a problem-dependent constant (large enough to allow the algorithm to “climb out” of regions around local minimizers that are not global minimizers). See also [57] for an analysis of a generalized version of simulated annealing.

Simulated annealing is often also used in combinatorial optimization, where the feasible set is finite (but typically large). An example of such a problem is the celebrated *traveling salesperson problem*. In the most basic form of this problem, we are given a number of cities and the cost of traveling from any city to any other city. The optimization problem is to find the cheapest round-trip route, starting from a given city, that visits every other city exactly once. For a description of how to apply simulated annealing to the traveling salesperson problem, see [67, p. 183].

#### 14.4 Particle Swarm Optimization

*Particle swarm optimization (PSO)* is a randomized search technique presented by James Kennedy (a social psychologist) and Russell C. Eberhart (an engineer) in 1995 [73]. This optimization method is inspired by social interaction principles. The PSO algorithm differs from the randomized search methods discussed in Section 14.3 in one key way: Instead of updating a single candidate solution  $\mathbf{x}^{(k)}$  at each iteration, we update a *population* (set) of candidate solutions, called a *swarm*. Each candidate solution in the swarm is called a *particle*. We think of a swarm as an apparently disorganized population of moving individuals that tend to cluster together while each individual seems to be moving in a random direction. (This description was adapted from a presentation by R. C. Eberhart.) The PSO algorithm aims to mimic the social behavior of animals and insects, such as a swarm of bees, a flock of birds, or a herd of wildebeest.

Suppose that we wish to minimize an objective function over  $\mathbb{R}^n$ . In the PSO algorithm, we start with an initial randomly generated population of points in  $\mathbb{R}^n$ . Associated with each point in the population is a velocity vector. We think of each point as the position of a particle, moving with an associated velocity. We then evaluate the objective function at each point in the population. Based on this evaluation, we create a new population of

points together with a new set of velocities. The creation of points in the new population, and their velocities, involve certain operations on points and velocities of the particles in the preceding population, described later.

Each particle keeps track of its *best-so-far position*—this is the best position it has visited so far (with respect to the value of the objective function). We will call this particle-dependent best-so-far position a *personal best* ( $pbest$ ). In contrast, the overall best-so-far position (best among *all* the positions encountered so far by the entire population) is called a *global best* ( $gbest$ ).

The particles “interact” with each other by updating their velocities according to their individual personal best as well as the global best. In the  $gbest$  version of the PSO algorithm, presented below, the velocity of each particle is changed, at each time step, toward a combination of its  $pbest$  and the  $gbest$  locations. The velocity is weighted by a random term, with separate random numbers being generated for velocities toward  $pbest$  and  $gbest$  locations. Thus, the particles are drawn both to their own personal best positions as well as to the best position of the entire swarm. As usual, typical stopping criteria of the algorithm consist of reaching a certain number of iterations, or reaching a certain objective function value.

### Basic PSO Algorithm

We now present a simple version of the  $gbest$  version of the PSO algorithm, where at each time step the velocity of each particle is changed toward its  $pbest$  and the  $gbest$  locations. Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be the objective function that we wish to minimize. Let  $d$  be the population size, and index the particles in the swarm by  $i = 1, \dots, d$ . Denote the position of particle  $i$  by  $\mathbf{x}_i \in \mathbb{R}^n$  and its velocity by  $\mathbf{v}_i \in \mathbb{R}^n$ . Let  $\mathbf{p}_i$  be the  $pbest$  of particle  $i$  and  $\mathbf{g}$  the  $gbest$ .

It is convenient to introduce the *Hadamard product* (or *Schur product*) operator, denoted by  $\circ$ : If  $\mathbf{A}$  and  $\mathbf{B}$  are matrices with the same dimension, then  $\mathbf{A} \circ \mathbf{B}$  is a matrix of the same dimension as  $\mathbf{A}$  (or  $\mathbf{B}$ ) resulting from entry-by-entry multiplication of  $\mathbf{A}$  and  $\mathbf{B}$ . This operation is denoted in MATLAB by “ $.*$ ” (the dot before an operator indicates entry-by-entry operations). Thus, if  $\mathbf{A}$  and  $\mathbf{B}$  have the same dimension, then  $\mathbf{A}.*\mathbf{B}$  returns a matrix whose entries are simply the products of the corresponding individual entries of  $\mathbf{A}$  and  $\mathbf{B}$ . The PSO  $gbest$  algorithm uses three given constant real parameters,  $\omega$ ,  $c_1$ , and  $c_2$ , which we discuss after presenting the algorithm.

### PSO Gbest Algorithm

- Set  $k := 0$ . For  $i = 1, \dots, d$ , generate initial random positions  $\mathbf{x}_i^{(0)}$  and velocities  $\mathbf{v}_i^{(0)}$ , and set  $\mathbf{p}_i^{(0)} = \mathbf{x}_i^{(0)}$ . Set  $\mathbf{g}^{(0)} = \arg \min_{\mathbf{x} \in \{\mathbf{x}_1^{(0)}, \dots, \mathbf{x}_d^{(0)}\}} f(\mathbf{x})$ .

2. For  $i = 1, \dots, d$ , generate random  $n$ -vectors  $\mathbf{r}_i^{(k)}$  and  $\mathbf{s}_i^{(k)}$  with components uniformly in the interval  $(0, 1)$ , and set
 
$$\begin{aligned}\mathbf{v}_i^{(k+1)} &= \omega \mathbf{v}_i^{(k)} + c_1 \mathbf{r}_i^{(k)} \circ (\mathbf{p}_i^{(k)} - \mathbf{x}_i^{(k)}) + c_2 \mathbf{s}_i^{(k)} \circ (\mathbf{g}^{(k)} - \mathbf{x}_i^{(k)}), \\ \mathbf{x}_i^{(k+1)} &= \mathbf{x}_i^{(k)} + \mathbf{v}_i^{(k+1)}.\end{aligned}$$
3. For  $i = 1, \dots, d$ , if  $f(\mathbf{x}_i^{(k+1)}) < f(\mathbf{p}_i^{(k)})$ , then set  $\mathbf{p}_i^{(k+1)} = \mathbf{x}_i^{(k+1)}$ ; else, set  $\mathbf{p}_i^{(k+1)} = \mathbf{p}_i^{(k)}$ .
4. If there exists  $i \in \{1, \dots, d\}$  such that  $f(\mathbf{x}_i^{(k+1)}) < f(\mathbf{g}^{(k)})$ , then set  $\mathbf{g}^{(k+1)} = \mathbf{x}_i^{(k+1)}$ ; else, set  $\mathbf{g}^{(k+1)} = \mathbf{g}^{(k)}$ .
5. If stopping criterion satisfied, then stop.
6. Set  $k := k + 1$ , go to step 2.

In the algorithm, the parameter  $\omega$  is referred to as an *inertial constant*. Recommended values are slightly less than 1. The parameters  $c_1$  and  $c_2$  are constants that determine how much the particle is directed toward “good” positions. They represent a “cognitive” and a “social” component, respectively, in that they affect how much the particle’s personal best and the global best influence its movement. Recommended values are  $c_1, c_2 \approx 2$ .

## Variations

The PSO techniques have evolved since 1995. For example, recently Clerc [29] proposed a *constriction-factor* version of the algorithm, where the velocity is updated as

$$\mathbf{v}_i^{(k+1)} = \kappa \left( \mathbf{v}_i^{(k)} + c_1 \mathbf{r}_i^{(k)} \circ (\mathbf{p}_i^{(k)} - \mathbf{x}_i^{(k)}) + c_2 \mathbf{s}_i^{(k)} \circ (\mathbf{g}^{(k)} - \mathbf{x}_i^{(k)}) \right),$$

where the *constriction coefficient*  $\kappa$  is computed as

$$\kappa = \frac{2}{|2 - \phi - \sqrt{\phi^2 - 4\phi}|},$$

where  $\phi = c_1 + c_2$  and  $\phi > 4$ . For example, for  $\phi = 4.1$ , we have  $\kappa = 0.729$ . The role of the constriction coefficient is to speed up the convergence.

When using PSO in practice, one might wish to clamp the velocities to a certain maximum amount, say,  $v_{\max}$ . In other words, we replace each component  $v$  of the velocity vector by

$$\min \{v_{\max}, \max\{-v_{\max}, v\}\}.$$

For an up-to-date literature survey and other modifications and heuristics, we recommend the first part of the proceedings of the *8th International Conference on Adaptive and Natural Computing Algorithms*, held in April 2007 in

Warsaw, Poland [5]. In these proceedings, one can find a number of papers dealing with applications of PSO to multiobjective optimization problems, versions of PSO for constrained optimization problems, as well as “niching” versions designed to find multiple solutions, that is, applications of PSO to multimodal optimization problems. For a mathematical analysis of the PSO algorithm, see Clerc and Kennedy [30].

## 14.5 Genetic Algorithms

### Basic Description

A *genetic algorithm* is a randomized, population-based search technique that has its roots in the principles of genetics. The beginnings of genetic algorithms is credited to John Holland, who developed the basic ideas in the late 1960s and early 1970s. Since its conception, genetic algorithms have been used widely as a tool in computer programming and artificial intelligence (e.g., [61], [79], and [94]), optimization (e.g., [36], [67], and [127]), neural network training (e.g., [80]), and many other areas.

Suppose that we wish to solve an optimization problem of the form

$$\begin{aligned} & \text{maximize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega \end{aligned}$$

(notice that the problem is a maximization, which is more convenient for describing genetic algorithms). The underlying idea of genetic algorithms applied to this problem is as follows. We start with an initial set of points in  $\Omega$ , denoted  $P(0)$ , called the *initial population*. We then evaluate the objective function at points in  $P(0)$ . Based on this evaluation, we create a new set of points  $P(1)$ . The creation of  $P(1)$  involves certain operations on points in  $P(0)$ , called *crossover* and *mutation*, discussed later. We repeat the procedure iteratively, generating populations  $P(2), P(3), \dots$ , until an appropriate stopping criterion is reached. The purpose of the crossover and mutation operations is to create a new population with an average objective function value that is higher than that of the previous population. To summarize, the genetic algorithm iteratively performs the operations of crossover and mutation on each population to produce a new population until a chosen stopping criterion is met.

The terminology used in describing genetic algorithms is adopted from genetics. To proceed with describing the details of the algorithm, we need the additional ideas and terms described below.

**Chromosomes and Representation Schemes** First, we point out that, in fact, genetic algorithms do not work directly with points in the set  $\Omega$ , but rather, with an *encoding* of the points in  $\Omega$ . Specifically, we need first to map  $\Omega$  onto a set consisting of strings of symbols, all of equal length. These strings are

called *chromosomes*. Each chromosome consists of elements from a chosen set of symbols, called the *alphabet*. For example, a common alphabet is the set  $\{0, 1\}$ , in which case the chromosomes are simply binary strings. We denote by  $L$  the length of chromosomes (i.e., the number of symbols in the strings). To each chromosome there corresponds a value of the objective function, referred to as the *fitness* of the chromosome. For each chromosome  $\mathbf{x}$ , we write  $f(\mathbf{x})$  for its fitness. Note that, for convenience, we use  $f$  to denote both the original objective function and the fitness measure on the set of chromosomes. We assume that  $f$  is a nonnegative function.

The choice of chromosome length, alphabet, and encoding (i.e., the mapping from  $\Omega$  onto the set of chromosomes) is called the *representation scheme* for the problem. Identification of an appropriate representation scheme is the first step in using genetic algorithms to solve a given optimization problem.

Once a suitable representation scheme has been chosen, the next phase is to initialize the first population  $P(0)$  of chromosomes. This is usually done by a random selection of a set of chromosomes. After we form the initial population of chromosomes, we then apply the operations of crossover and mutation on the population. During each iteration  $k$  of the process, we evaluate the fitness  $f(\mathbf{x}^{(k)})$  of each member  $\mathbf{x}^{(k)}$  of the population  $P(k)$ . After the fitness of the entire population has been evaluated, we form a new population  $P(k+1)$  in two stages.

**Selection and Evolution** In the first stage we apply an operation called *selection*, where we form a set  $M(k)$  with the same number of elements as  $P(k)$ . This number is called the *population size*, which we denote by  $N$ . The set  $M(k)$ , called the *mating pool*, is formed from  $P(k)$  using a random procedure as follows: Each point  $\mathbf{m}^{(k)}$  in  $M(k)$  is equal to  $\mathbf{x}^{(k)}$  in  $P(k)$  with probability

$$\frac{f(\mathbf{x}^{(k)})}{F(k)},$$

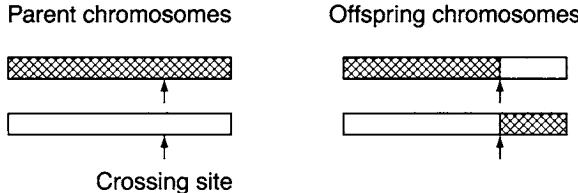
where

$$F(k) = \sum f(\mathbf{x}_i^{(k)})$$

and the sum is taken over the whole of  $P(k)$ . In other words, we select chromosomes into the mating pool with probabilities proportional to their fitness.

The selection scheme above is also called the *roulette-wheel* scheme, for the following reason. Imagine a roulette wheel in which each slot is assigned to a chromosome in  $P(k)$ ; some chromosomes may be assigned multiple slots. The number of slots associated with each chromosome is in proportion to its fitness. We then spin the roulette wheel and select [for inclusion in  $M(k)$ ] the chromosome on whose slot the ball comes to rest. This procedure is repeated  $N$  times, so that the mating pool  $M(k)$  contains  $N$  chromosomes.

An alternative selection scheme is the *tournament scheme*, which proceeds as follows. First, we select a pair of chromosomes at random from  $P(k)$ . We



**Figure 14.7** Illustration of basic crossover operation.

then compare the fitness values of these two chromosomes, and place the fitter of the two into  $M(k)$ . We repeat this operation until the mating pool  $M(k)$  contains  $N$  chromosomes.

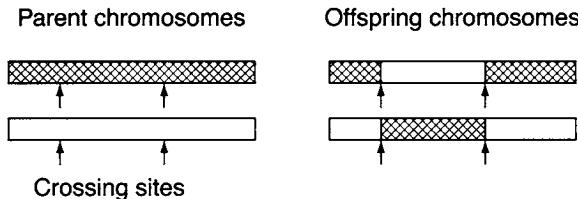
The second stage is called *evolution*: in this stage, we apply the crossover and mutation operations. The *crossover operation* takes a pair of chromosomes, called the *parents*, and gives a pair of *offspring chromosomes*. The operation involves exchanging substrings of the two parent chromosomes, described below. Pairs of parents for crossover are chosen from the mating pool randomly, such that the probability that a chromosome is chosen for crossover is  $p_c$ . We assume that whether or not a given chromosome is chosen is independent of whether or not any other chromosome is chosen for crossover.

We can pick parents for crossover in several ways. For example, we may randomly choose two chromosomes from the mating pool as parents. In this case if  $N$  is the number of chromosomes in the mating pool, then  $p_c = 2/N$ . Similarly, if we randomly pick  $2k$  chromosomes from the mating pool (where  $k < N/2$ ), forming  $k$  pairs of parents, we have  $p_c = 2k/N$ . In the two examples above, the number of pairs of parents is fixed and the value of  $p_c$  is dependent on this number. Yet another way of choosing parents is as follows: Given a value of  $p_c$ , we pick a random number of pairs of parents such that the average number of pairs is  $p_cN/2$ .

Once the parents for crossover have been determined, we apply the crossover operation to the parents. There are many types of possible crossover operations. The simplest crossover operation is the *one-point crossover*. In this operation, we first choose a number randomly between 1 and  $L - 1$  according to a uniform distribution, where  $L$  is the length of chromosomes. We refer to this number as the *crossing site*. Crossover then involves exchanging substrings of the parents to the left of the crossing site, as illustrated in Figure 14.7 and in the following example.

**Example 14.1** Suppose that we have chromosomes of length  $L = 6$  over the binary alphabet  $\{0, 1\}$ . Consider the pair of parents 000000 and 111111. Suppose that the crossing site is 4. Then, the crossover operation applied to the parent chromosomes yields the two offspring 000011 and 111100. ■

We can also have crossover operations with multiple crossing sites, as illustrated in Figure 14.8 and in the following example.



**Figure 14.8** Illustration of two-point crossover operation.

**Example 14.2** Consider two chromosomes, 000000000 and 111111111, of length  $L = 9$ . Suppose that we have two crossing sites, at 3 and 7. Then, the crossover operation applied to the parent chromosomes above yields the two offspring 000111100 and 111000011. ■

After the crossover operation, we replace the parents in the mating pool by their offspring. The mating pool has therefore been modified but maintains the same number of elements.

Next, we apply the *mutation operation*, which takes each chromosome from the mating pool and randomly changes each symbol of the chromosome with a given probability  $p_m$ . In the case of the binary alphabet, this change corresponds to complementing the corresponding bits; that is, we replace each bit with probability  $p_m$  from 0 to 1, or vice versa. If the alphabet contains more than two symbols, then the change involves randomly substituting the symbol with another symbol from the alphabet. Typically, the value of  $p_m$  is very small (e.g., 0.01), so that only a few chromosomes will undergo a change due to mutation, and of those that are affected, only a few of the symbols are modified. Therefore, the mutation operation plays only a minor role in the genetic algorithm relative to the crossover operation.

After applying the crossover and mutation operations to the mating pool  $M(k)$ , we obtain the new population  $P(k+1)$ . We then repeat the procedure of evaluation, selection, and evolution, iteratively. We summarize the genetic algorithm as follows.

### Genetic Algorithm

1. Set  $k := 0$ . Generate an initial population  $P(0)$ .
2. Evaluate  $P(k)$ .
3. If the stopping criterion is satisfied, then stop.
4. Select  $M(k)$  from  $P(k)$ .
5. Evolve  $M(k)$  to form  $P(k + 1)$ .
6. Set  $k := k + 1$ , go to step 2.

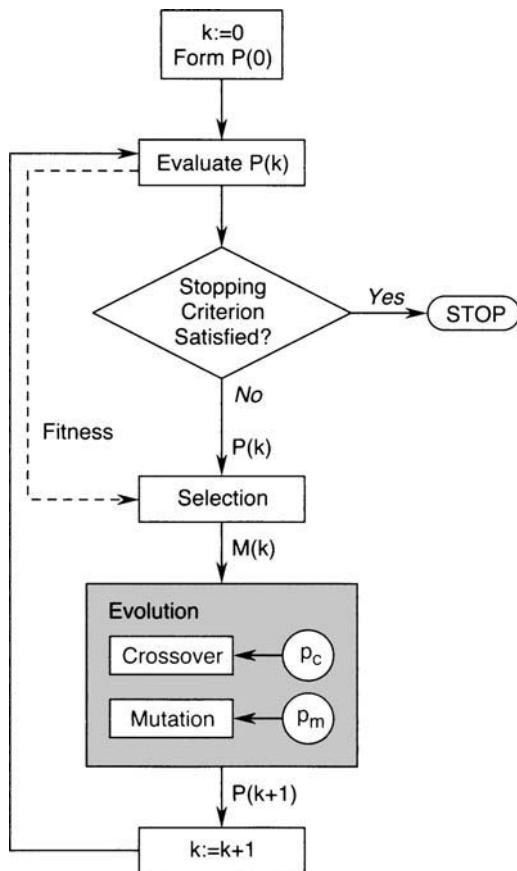


Figure 14.9 Flowchart for the genetic algorithm.

A flowchart illustrating this algorithm is shown in Figure 14.9.

During execution of the genetic algorithm, we keep track of the *best-so-far* chromosome, that is, the chromosome with the highest fitness of all the chromosomes evaluated. After each iteration, the best-so-far chromosome serves as the candidate for the solution to the original problem. Indeed, we may even copy the best-so-far chromosome into each new population, a practice referred to as *elitism*. The elitist strategy may result in domination of the population by “superchromosomes.” However, practical experience suggests that elitism often improves the performance of the algorithm.

The stopping criterion can be implemented in a number of ways. For example, a simple stopping criterion is to stop after a prespecified number of iterations. Another possible criterion is to stop when the fitness for the best-so-far chromosome does not change significantly from iteration to iteration.

The genetic algorithm differs from the algorithms discussed in previous chapters in several respects. First, it does not use derivatives of the objective function (like the other methods in this chapter). Second, it uses operations that are random within each iteration (like the other randomized search methods). Third, it searches from a set of points rather than a single point at each iteration (like the PSO algorithm). Fourth, it works with an encoding of the feasible set rather with than the set itself.

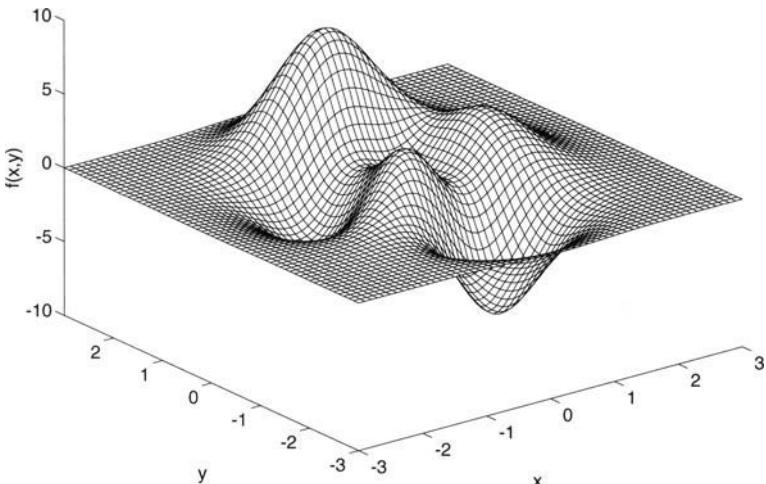
We illustrate an application of the genetic algorithm to an optimization problem in the following example.

**Example 14.3** Consider the MATLAB “peaks” function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by

$$f(x, y) = 3(1 - x)^2 e^{-x^2 - (y+1)^2} - 10 \left( \frac{x}{5} - x^3 - y^5 \right) e^{-x^2 - y^2} - \frac{e^{-(x+1)^2 - y^2}}{3}$$

(see also [67, pp. 178–180] for an example involving the same function). We wish to maximize  $f$  over the set  $\Omega = \{[x, y]^\top \in \mathbb{R}^2 : -3 \leq x, y \leq 3\}$ . A plot of the objective function  $f$  over the feasible set  $\Omega$  is shown in Figure 14.10. Using the MATLAB function `fminunc` (from the Optimization Toolbox), we found the optimal point to be  $[-0.0093, 1.5814]^\top$ , with objective function value 8.1062.

To apply the genetic algorithm to solve the optimization problem above, we use a simple binary representation scheme with length  $L = 32$ , where the first 16 bits of each chromosome encode the  $x$  component, whereas the remaining 16 bits encode the  $y$  component. Recall that  $x$  and  $y$  take values in the interval  $[-3, 3]$ . We first map the interval  $[-3, 3]$  onto the interval  $[0, 2^{16} - 1]$ , via a simple translation and scaling. The integers in the interval  $[0, 2^{16} - 1]$  are then expressed as binary 16-bit strings. This defines the encoding of each component  $x$  and  $y$ . The chromosome is obtained by juxtaposing the two 8-bit strings. For example, the point  $[x, y]^\top = [-1, 3]^\top$  is encoded as (see



**Figure 14.10** Plot of  $f$  for Example 14.3.

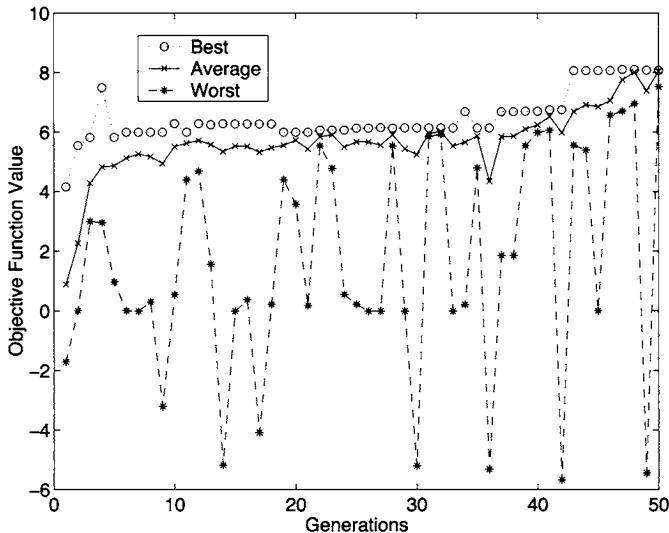
Exercise 14.4 for a simple algorithm for converting from decimal into binary)

$$\underbrace{01010101010101}_{\text{encoded } x = -1} \underbrace{1111111111111111}_{\text{encoded } y = 3}.$$

Using a population size of 20, we apply 50 iterations of the genetic algorithm on the problem above. We used parameter values of  $p_c = 0.75$  and  $p_m = 0.0075$ . Figure 14.11 shows plots of the best, average, and worst objective function values in the population for every iteration (generation) of the algorithm. The best-so-far solution obtained at the end of the 50 iterations is  $[0.0615, 1.5827]^\top$ , with objective function value 8.1013. Note that this solution and objective function value are very close to those obtained using MATLAB. ■

### Analysis of Genetic Algorithms

In this section we use heuristic arguments to describe why genetic algorithms work. As pointed out before, the genetic algorithm was motivated by ideas from natural genetics [61]. Specifically, the notion of “survival of the fittest” plays a central role. The mechanisms used in the genetic algorithm mimic this principle. We start with a population of chromosomes, and selectively pick the fittest ones for reproduction. From these selected chromosomes, we form the new generation by combining information encoded in them. In this way,



**Figure 14.11** The best, average, and worst objective function values in the population for every iteration (generation) of the genetic algorithm in Example 14.3.

the goal is to ensure that the fittest members of the population survive and their information content is preserved and combined to produce even better offspring.

To further analyze the genetic algorithm in a more quantitative fashion, we need to define a few terms. For convenience, we only consider chromosomes over the binary alphabet. We introduce the notion of a *schema* (plural: *schemata*) as a set of chromosomes with certain common features. Specifically, a schema is a set of chromosomes that contain 1s and 0s in particular locations. We represent a schema using a string notation over an extended alphabet  $\{0, 1, *\}$ . For example, the notation  $1 * 01$  represents the schema

$$1 * 01 = \{1001, 1101\},$$

and the notation  $0 * 101*$  represents the schema

$$0 * 101* = \{001010, 001011, 011010, 011011\}.$$

In the schema notation, the numbers 0 and 1 denote the fixed binary values in the chromosomes that belong to the schema. The symbol \*, meaning “don’t care,” matches either 0 or 1 at the positions it occupies. Thus, a schema describes a set of chromosomes that have certain specified similarities. A chromosome belongs to a particular schema if for all positions  $j = 1, \dots, L$  the symbol found in the  $j$ th position of the chromosome matches the symbol found in the  $j$ th position of the schema, with the understanding that any symbol matches \*. Note that if a schema has  $r$  “don’t care” symbols, then it

contains  $2^r$  chromosomes. Moreover, any chromosome of length  $L$  belongs to  $2^L$  schemata.

Given a schema that represents good solutions to our optimization problem, we would like the number of matching chromosomes in the population  $P(k)$  to grow as  $k$  increases. This growth is affected by several factors, which we analyze in the following discussion. We assume throughout that we are using the roulette-wheel selection method.

The first key idea in explaining why the genetic algorithm works is the observation that if a schema has chromosomes with better-than-average fitness, then the expected (mean) number of chromosomes matching this schema in the mating pool  $M(k)$  is larger than the number of chromosomes matching this schema in the population  $P(k)$ . To quantify this assertion, we need some additional notation. Let  $H$  be a given schema, and let  $e(H, k)$  be the number of chromosomes in  $P(k)$  that match  $H$ ; that is,  $e(H, k)$  is the number of elements in the set  $P(k) \cap H$ . Let  $f(H, k)$  be the average fitness of chromosomes in  $P(k)$  that match schema  $H$ . This means that if  $H \cap P(k) = \{\mathbf{x}_1, \dots, \mathbf{x}_{e(H, k)}\}$ , then

$$f(H, k) = \frac{f(\mathbf{x}_1) + \dots + f(\mathbf{x}_{e(H, k)})}{e(H, k)}.$$

Let  $N$  be the number of chromosomes in the population and  $F(k)$  be the sum of the fitness values of chromosomes in  $P(k)$ , as before. Denote by  $\bar{F}(k)$  the average fitness of chromosomes in the population; that is,

$$\begin{aligned}\bar{F}(k) &= \frac{F(k)}{N} \\ &= \frac{1}{N} \sum f(\mathbf{x}_i^{(k)}).\end{aligned}$$

Finally, let  $m(H, k)$  be the number of chromosomes in  $M(k)$  that match  $H$ , in other words, the number of elements in the set  $M(k) \cap H$ .

**Lemma 14.1** *Let  $H$  be a given schema and  $\mathcal{M}(H, k)$  be the expected value of  $m(H, k)$  given  $P(k)$ . Then,*

$$\mathcal{M}(H, k) = \frac{f(H, k)}{\bar{F}(k)} e(H, k).$$

□

*Proof.* Let  $P(k) \cap H = \{\mathbf{x}_1, \dots, \mathbf{x}_{e(H, k)}\}$ . In the remainder of the proof, the term *expected* should be taken to mean “expected, given  $P(k)$ .” For each element  $\mathbf{m}^{(k)} \in M(k)$  and each  $i = 1, \dots, e(H, k)$ , the probability that  $\mathbf{m}^{(k)} = \mathbf{x}_i$  is given by  $f(\mathbf{x}_i)/F(k)$ . Thus, the expected number of chromosomes in  $M(k)$  equal to  $\mathbf{x}_i$  is

$$N \frac{f(\mathbf{x}_i)}{F(k)} = \frac{f(\mathbf{x}_i)}{\bar{F}(k)}.$$

Hence, the expected number of chromosomes in  $P(k) \cap H$  that are selected into  $M(k)$  is

$$\sum_{i=1}^{e(H,k)} \frac{f(\mathbf{x}_i)}{\bar{F}(k)} = e(H,k) \frac{\sum_{i=1}^{e(H,k)} f(\mathbf{x}_i)}{e(H,k)} \frac{1}{\bar{F}(k)} = \frac{f(H,k)}{\bar{F}(k)} e(H,k).$$

Because any chromosome in  $M(k)$  is also a chromosome in  $P(k)$ , the chromosomes in  $M(k) \cap H$  are simply those in  $P(k) \cap H$  that are selected into  $M(k)$ . Hence,

$$\mathcal{M}(H, k) = \frac{f(H, k)}{\bar{F}(k)} e(H, k).$$

■

Lemma 14.1 quantifies our assertion that if a schema  $H$  has chromosomes with better than average fitness [i.e.,  $f(H, k)/\bar{F}(k) > 1$ ], then the expected number of chromosomes matching  $H$  in the mating pool  $M(k)$  is larger than the number of chromosomes matching  $H$  in the population  $P(k)$ .

We now analyze the effect of the evolution operations on the chromosomes in the mating pool. For this, we need to introduce two parameters that are useful in the characterization of a schema: order and length. The *order*  $o(S)$  of a schema  $S$  is the number of fixed symbols (non\* symbols) in its representation (the notation  $o(S)$  is standard in the literature on genetic algorithms, and should not be confused with the “little-oh” symbol defined in Section 5.6). If the *length* of chromosomes in  $S$  is  $L$ , then  $o(S)$  is  $L$  minus the number of \* symbols in  $S$ . For example,

$$o(1 * 01) = 4 - 1 = 3,$$

whereas

$$o(0 * 1 * 01) = 6 - 2 = 4.$$

The *length*  $l(S)$  of a schema  $S$  is the distance between the first and last fixed symbols (i.e., the difference between the positions of the rightmost fixed symbol and the leftmost fixed symbol). For example,

$$l(1 * 01) = 4 - 1 = 3,$$

$$l(0 * 101*) = 5 - 1 = 4,$$

$$l(**1*) = 0.$$

Note that for a schema  $S$  with chromosomes of length  $L$ , the order  $o(S)$  is a number between 0 and  $L$  and the length  $l(S)$  is a number between 0 in  $L - 1$ . The order of a schema with all \* symbols is 0; its length is also 0. The order of a schema containing only a single element (i.e., its representation has no \* symbols) is  $L$  [e.g.,  $o(1011) = 4 - 0 = 4$ ]. The length of a schema with fixed symbols in its first and last positions is  $L - 1$  [e.g.,  $l(0 ** 1) = 4 - 1 = 3$ ].

We first consider the effect of the crossover operation on the mating pool. The basic observation in the following lemma is that given a chromosome in  $M(k) \cap H$ , the probability that it leaves  $H$  after crossover is bounded above by a quantity that is proportional to  $p_c$  and  $l(H)$ .

**Lemma 14.2** *Given a chromosome in  $M(k) \cap H$ , the probability that it is chosen for crossover and neither of its offspring is in  $H$  is bounded above by*

$$p_c \frac{l(H)}{L - 1}.$$

□

*Proof.* Consider a given chromosome in  $M(k) \cap H$ . The probability that it is chosen for crossover is  $p_c$ . If neither of its offspring is in  $H$ , then the crossover point must be between the corresponding first and last fixed symbols of  $H$ . The probability of this is  $l(H)/(L - 1)$ . Hence, the probability that the given chromosome is chosen for crossover and neither of its offspring is in  $H$  is bounded above by

$$p_c \frac{l(H)}{L - 1}.$$

■

From Lemma 14.2 we conclude that given a chromosome in  $M(k) \cap H$ , the probability either that it is not selected for crossover or that at least one of its offspring is in  $H$  after the crossover operation, is bounded below by

$$1 - p_c \frac{l(H)}{L - 1}.$$

Note that if a chromosome in  $H$  is chosen for crossover and the other parent chromosome is also in  $H$ , then both offspring are automatically in  $H$  (see Exercise 14.5). Hence, for each chromosome in  $M(k) \cap H$ , there is a certain probability that it will result in an associated chromosome in  $H$  (either itself or one of its offspring) after going through crossover (including selection for crossover) and that probability is bounded below by the foregoing expression.

We next consider the effect of the mutation operation on the mating pool  $M(k)$ .

**Lemma 14.3** *Given a chromosome in  $M(k) \cap H$ , the probability that it remains in  $H$  after the mutation operation is given by*

$$(1 - p_m)^{o(H)}.$$

□

*Proof.* Given a chromosome in  $M(k) \cap H$ , it remains in  $H$  after the mutation operation if and only if none of the symbols in this chromosome that correspond to fixed symbols in  $H$  are changed by the mutation operation. The probability of this event is  $(1 - p_m)^{o(H)}$ . ■

Note that if  $p_m$  is small, the expression  $(1 - p_m)^{o(H)}$  above is approximately equal to

$$1 - p_m o(H).$$

The following theorem combines the results of the preceding lemmas.

**Theorem 14.1** *Let  $H$  be a given schema and  $\mathcal{E}(H, k + 1)$  be the expected value of  $e(H, k + 1)$  given  $P(k)$ . Then,*

$$\mathcal{E}(H, k + 1) \geq \left(1 - p_c \frac{l(H)}{L - 1}\right) (1 - p_m)^{o(H)} \frac{f(H, k)}{\bar{F}(k)} e(H, k).$$

□

*Proof.* Consider a given chromosome in  $M(k) \cap H$ . If, after the evolution operations, it has a resulting chromosome that is in  $H$ , then that chromosome is in  $P(k + 1) \cap H$ . By Lemmas 14.2 and 14.3, the probability of this event is bounded below by

$$\left(1 - p_c \frac{l(H)}{L - 1}\right) (1 - p_m)^{o(H)}.$$

Therefore, because each chromosome in  $M(k) \cap H$  results in a chromosome in  $P(k + 1) \cap H$  with a probability bounded below by the expression above, the expected value of  $e(H, k + 1)$  given  $M(k)$  is bounded below by

$$\left(1 - p_c \frac{l(H)}{L - 1}\right) (1 - p_m)^{o(H)} m(H, k).$$

Taking the expectation given  $P(k)$ , we get

$$\mathcal{E}(H, k + 1) \geq \left(1 - p_c \frac{l(H)}{L - 1}\right) (1 - p_m)^{o(H)} \mathcal{M}(H, k).$$

Finally, using Lemma 14.1, we arrive at the desired result. ■

Theorem 14.1 indicates how the number of chromosomes in a given schema changes from one population to the next. Three factors influence this change, reflected by the three terms on the right-hand side of inequality in Theorem 14.1:  $1 - p_c l(H)/(L - 1)$ ,  $(1 - p_m)^{o(H)}$ , and  $f(H, k)/\bar{F}(k)$ . Note that the larger the values of these terms, the higher the expected number of matches of the schema  $H$  in the next population. The effect of each term is summarized as follows:

- The term  $f(H, k)/\bar{F}(k)$  reflects the role of average fitness of the given schema  $H$ —the higher the average fitness, the higher the expected number of matches in the next population.
- The term  $1 - p_c l(H)/(L - 1)$  reflects the effect of crossover—the smaller the term  $p_c l(H)/(L - 1)$ , the higher the expected number of matches in the next population.

- The term  $(1-p_m)^{o(H)}$  reflects the effect of mutation—the larger the term, the higher the expected number of matches in the next population.

In summary, we see that a schema that is short, low order, and has above-average fitness will have on average an increasing number of its representatives in the population from iteration to iteration. Observe that the encoding is relevant to the performance of the algorithm. Specifically, a good encoding is one that results in high-fitness schemata having small lengths and orders.

### Real-Number Genetic Algorithms

The genetic algorithms described thus far operate on binary strings, representing elements of the feasible set  $\Omega$ . (For this reason, genetic algorithms are also suitably applied to combinatorial optimization problems, where  $\Omega$  is not  $\mathbb{R}^n$  but some discrete set.) Binary encodings allow us to use the schema theory, described in the preceding section, to analyze genetic algorithms. However, there are some disadvantages to operating on binary strings. To see this, let  $\mathbf{g} : \{0, 1\}^L \rightarrow \Omega$  represent the binary “decoding” function; that is, if  $\mathbf{x}$  is a binary chromosome,  $\mathbf{g}(\mathbf{x}) \in \Omega$  is the point in the feasible set  $\Omega \subset \mathbb{R}^n$  whose encoding is  $\mathbf{x}$ . Therefore, the objective function being maximized by the genetic algorithm is not  $f$  itself but rather the composition of  $f$  and the decoding function  $\mathbf{g}$ . In other words, the optimization problem being solved by the genetic algorithm is

$$\begin{aligned} &\text{maximize} && f(\mathbf{g}(\mathbf{x})) \\ &\text{subject to} && \mathbf{x} \in \{\mathbf{y} \in \{0, 1\}^L : \mathbf{g}(\mathbf{y}) \in \Omega\}. \end{aligned}$$

This optimization problem may be more complex than the original optimization problem. For example, it may have extra maximizers, making the search for a global maximizer more difficult.

The above motivates a consideration of genetic algorithms that operate directly on the original optimization problem. In other words, we wish to implement a genetic algorithm that operates directly on  $\mathbb{R}^n$ . The steps of this algorithm will be the same as before (see Figure 14.9), except that the elements of the population are points in the feasible set  $\Omega$  rather than binary strings. We will need to define appropriate crossover and mutation operations for this case.

For crossover, we have several options. The simplest is to use averaging: For a pair of parents  $\mathbf{x}$  and  $\mathbf{y}$ , the offspring is  $\mathbf{z} = (\mathbf{x} + \mathbf{y})/2$  (this type of crossover operation is used, e.g., in [103]). This offspring can then replace one of the parents. Alternatively, we may produce two offspring as follows:  $\mathbf{z}_1 = (\mathbf{x} + \mathbf{y})/2 + \mathbf{w}_1$  and  $\mathbf{z}_2 = (\mathbf{x} + \mathbf{y})/2 + \mathbf{w}_2$ , where  $\mathbf{w}_1$  and  $\mathbf{w}_2$  are two randomly generated vectors (with zero mean). If either offspring lies outside  $\Omega$ , we have to bring the offspring back into  $\Omega$ , using, for example, a projection (see Section 23.2). A third option for crossover is to take random convex combinations of the parents. Specifically, given a pair of parents  $\mathbf{x}$  and  $\mathbf{y}$ ,

we generate a random number  $\alpha \in (0, 1)$  and then produce two offspring  $\mathbf{z}_1 = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$  and  $\mathbf{z}_2 = (1 - \alpha)\mathbf{x} + \alpha\mathbf{y}$ . This method of crossover ensures that the offspring are always in the feasible set, provided that the feasible set is convex. A fourth option is to perturb the two points above by some random amount:  $\mathbf{z}_1 = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} + \mathbf{w}_1$  and  $\mathbf{z}_2 = (1 - \alpha)\mathbf{x} + \alpha\mathbf{y} + \mathbf{w}_2$ , where  $\mathbf{w}_1$  and  $\mathbf{w}_2$  are two randomly generated vectors (with zero mean). In this case we have to check for feasibility of the offspring and use projections if needed.

For mutation, a simple implementation is to add a random vector to the chromosome. Specifically, given a chromosome  $\mathbf{x}$ , we produce its mutation as  $\mathbf{x}' = \mathbf{x} + \mathbf{w}$ , where  $\mathbf{w}$  is a random vector with zero mean. This mutation operation is also called a *real number creep* (see, e.g., [103]). As before, we have to ensure that the mutated chromosome is feasible. If not, we may use a projection. An alternative method for mutation is to replace the chosen chromosome with a random convex combination of the chromosome with a random point in the feasible set; that is, we generate a random number  $\alpha \in (0, 1)$  and a random point  $\mathbf{w} \in \Omega$ , and set  $\mathbf{x}' = \alpha\mathbf{x} + (1 - \alpha)\mathbf{w}$ . Provided that the feasible set is convex, the mutated chromosome will always be feasible.

**Example 14.4** Consider again the function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  from Example 14.3. We apply a real-number genetic algorithm to find a maximizer of  $f$  using a crossover operation of the fourth type described above and a mutation operation of the second type above. With a population size of 20, we apply 50 iterations of the genetic algorithm. As before, we used parameter values of  $p_c = 0.75$  and  $p_m = 0.0075$ . Figure 14.12 shows plots of the best, average, and worst objective function values in the population for every iteration (generation) of the algorithm. The best-so-far solution obtained at the end of the 50 iterations is  $[-0.0096, 1.5845]^\top$ , with objective function value 8.1061, which is close to the result of Example 14.3. ■

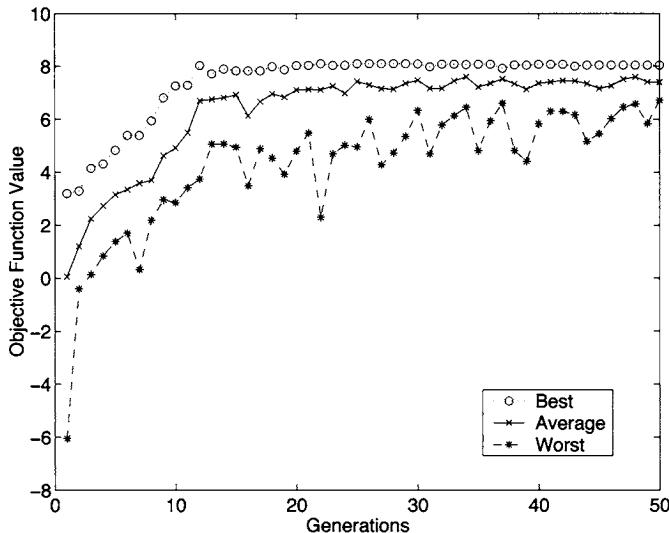
## EXERCISES

**14.1** Write a MATLAB program to implement the Nelder-Mead algorithm applied to minimizing the function

$$f(x_1, x_2) = (x_2 - x_1)^4 + 12x_1x_2 - x_1 + x_2 - 3$$

on  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1, x_2 \in [-1, 1]\}$ . Locate the iteration points on the level sets of  $f$ . Connect the successive points with lines to show clearly the progression of the optimization process. Test your program with two starting points:

$$\mathbf{x}^{(0)} = \begin{bmatrix} 0.55 \\ 0.7 \end{bmatrix} \quad \text{and} \quad \mathbf{x}^{(0)} = \begin{bmatrix} -0.9 \\ -0.5 \end{bmatrix}.$$



**Figure 14.12** The best, average, and worst objective function values in the population for every iteration (generation) of the real-number genetic algorithm in Example 14.4.

**14.2** Write MATLAB programs to implement naive random search and simulated annealing. Use the neighborhood

$$N(\mathbf{x}^{(k)}) = \{\mathbf{x} : x_i^{(k)} - \alpha \leq x_i \leq x_i^{(k)} + \alpha\},$$

where  $\alpha > 0$  is prespecified, and pick  $\mathbf{z}^{(k)}$  to be uniformly distributed on  $N(\mathbf{x}^{(k)})$ . Test both algorithms on maximizing the MATLAB “peaks” function given in Example 14.3. Observe the effect of varying  $\alpha$ .

**14.3** Write a MATLAB program to implement a particle swarm optimization algorithm. Test your implementation on maximizing the MATLAB “peaks” function given in Example 14.3.

**14.4** This problem has four parts and is related to binary encoding for genetic algorithms.

- a. Let  $(I)_{10}$  be the decimal representation for a given integer, and let  $a_m a_{m-1} \cdots a_0$  be its binary representation; that is, each  $a_i$  is either 0 or 1, and

$$(I)_{10} = a_m 2^m + a_{m-1} 2^{m-1} + \cdots + a_1 2^1 + a_0 2^0.$$

Verify that the following is true:

$$(I)_{10} = (((\cdots (((a_m 2 + a_{m-1}) 2 + a_{m-2}) 2 + a_{m-3}) \cdots ) 2 + a_1) 2 + a_0).$$

- b. The second expression in part a suggests a simple algorithm for converting from decimal representation to equivalent binary representation, as follows. Dividing both sides of the expression in part a by 2, the remainder is  $a_0$ . Subsequent divisions by 2 yield the remaining bits  $a_1, a_2, \dots, a_m$  as remainders.

Use this algorithm to find the binary representation of the integer  $(I)_{10} = 1995$ .

- c. Let  $(F)_{10}$  be the decimal representation for a given number in  $[0, 1]$ , and let  $0.a_{-1}a_{-2}\dots$  be its binary representation, that is,

$$(F)_{10} = a_{-1}2^{-1} + a_{-2}2^{-2} + \dots$$

If this expression is multiplied by 2, the integer part of the product is  $a_{-1}$ . Subsequent multiplications yield the remaining bits  $a_{-2}, a_{-3}, \dots$ . As in part b, the above gives a simple algorithm for converting from a decimal fraction to its binary representations. Use this algorithm to find the binary representation of  $(F)_{10} = 0.7265625$ .

Note that we can combine the algorithms from parts b and c to convert an arbitrary positive decimal representation into its equivalent binary representation. Specifically, we apply the algorithms in parts b and c separately to the integer and fraction parts of the given decimal number, respectively.

- d. The procedure in part c may yield an infinitely long binary representation. If this is the case, then we need to determine the number of bits required to keep at least the same accuracy as the given decimal number. If we have a  $d$ -digit decimal fraction, then the number of bits  $b$  in the binary representation must satisfy  $2^{-b} \leq 10^{-d}$ , which yields  $b \geq 3.32d$ . Convert 19.95 to its equivalent binary representation with at least the same degree of accuracy (i.e., to two decimal places).

**14.5** Given two chromosomes in a schema  $H$ , suppose that we swap some (or all) of the symbols between them at corresponding positions. Show that the resulting two chromosomes are also in  $H$ . From this fact we conclude that given two chromosomes in  $H$ , both offspring after the crossover operation are also in  $H$ . In other words, the crossover operation preserves membership in  $H$ .

**14.6** Consider a two-point crossover scheme (see Example 14.2), described as follows. Given a pair of binary chromosomes of length  $L$ , we independently choose two random numbers, uniform over  $1, \dots, L - 1$ . We call the two numbers  $c_1$  and  $c_2$ , where  $c_1 \leq c_2$ . If  $c_1 = c_2$ , we do not swap any symbols (i.e., leave the two given parent chromosomes unchanged). If  $c_1 < c_2$ , we interchange the  $(c_1 + 1)$ th through  $c_2$ th bits in the given parent chromosomes.

Prove the analog of Lemma 14.2 for this case, given below.

**Lemma:** Given a chromosome in  $M(k) \cap H$ , the probability that it is chosen for crossover and neither of its offspring is in  $H$  is bounded above by

$$p_c \left( 1 - \left( 1 - \frac{l(H)}{L-1} \right)^2 \right).$$

□

*Hint:* Note that the two-point crossover operation is equivalent to a composition of two one-point crossover operations (i.e., doing two one-point crossover operations in succession).

**14.7** State and prove the analog of Lemma 14.2 for an  $n$ -point crossover operation.

*Hint:* See Exercise 14.6.

**14.8** Implement the roulette-wheel selection scheme using MATLAB.

*Hint:* Use the MATLAB functions `sum`, `cumsum`, and `find`.

**14.9** Implement the crossover operation (one-point) using the MATLAB, assuming that we are given two binary parent chromosomes.

**14.10** Implement the mutation operation using the MATLAB function `xor`, assuming that the chromosomes in the mating pool are binary vectors.

**14.11** Write a MATLAB program to implement a genetic algorithm using binary encoding. Test your implementation on the following functions:

a.  $f(x) = -15 \sin^2(2x) - (x-2)^2 + 160$ ,  $|x| \leq 10$ .

b.  $f(x, y) = 3(1-x)^2 e^{-x^2-(y+1)^2} - 10 \left( \frac{x}{5} - x^3 - y^5 \right) e^{-x^2-y^2} - \frac{e^{-(x+1)^2-y^2}}{3}$ ,  
 $|x|, |y| \leq 3$  (considered in Example 14.3).

**14.12** Write a MATLAB program to implement a real-number genetic algorithm. Test your implementation on the function  $f(\mathbf{x}) = x_1 \sin(x_1) + x_2 \sin(5x_2)$  with the constraint set  $\Omega = \{\mathbf{x} : 0 \leq x_1 \leq 10, 4 \leq x_2 \leq 6\}$ .



## **PART III**

---

# **LINEAR PROGRAMMING**

---



## CHAPTER 15

---

# INTRODUCTION TO LINEAR PROGRAMMING

---

### 15.1 Brief History of Linear Programming

The goal of linear programming is to determine the values of decision variables that maximize or minimize a linear objective function, where the decision variables are subject to linear constraints. A linear programming problem is a special case of a general constrained optimization problem. In the general setting, the goal is to find a point that minimizes the objective function and at the same time satisfies the constraints. We refer to any point that satisfies the constraints as a *feasible point*. In a linear programming problem, the objective function is linear, and the set of feasible points is determined by a set of linear equations and/or inequalities.

In this part we study methods for solving linear programming problems. Linear programming methods provide a way of choosing the best feasible point among the many possible feasible points. In general, the number of feasible points is infinitely large. However, as we shall see, the solution to a linear programming problem can be found by searching through a particular finite number of feasible points, known as *basic feasible solutions*. Therefore, in principle, we can solve a linear programming problem simply by comparing

the finite number of basic feasible solutions and finding one that minimizes or maximizes the objective function—we refer to this approach as the *brute-force approach*. For most practical decision problems, even this finite number of basic feasible solutions is so large that the method of choosing the best solution by comparing them to each other is impractical. To get a feel for the amount of computation needed in a brute-force approach, consider the following example. Suppose that we have a small factory with 20 different machines producing 20 different parts. Assume that any of the machines can produce any part. We also assume that the time for producing each part on each machine is known. The problem then is to assign a part to each machine so that the overall production time is minimized. We see that there are  $20!$  (20 factorial) possible assignments. The brute-force approach to solving this assignment problem would involve writing down all the possible assignments and then choosing the best one by comparing them. Suppose that we have at our disposal a computer that takes  $1 \mu\text{s}$  ( $10^{-6}$  second) to determine each assignment. Then, to find the best (optimal) assignment this computer would need 77,147 years (working 24 hours a day, 365 days a year) to find the best solution. An alternative approach to solving this problem is to use experienced planners to optimize this assignment problem. Such an approach relies on heuristics. Heuristics come close, but give suboptimal answers. Heuristics that do reasonably well, with an error of, say, 10%, may still not be good enough. For example, in a business that operates on large volumes and a small profit margin, a 10% error could mean the difference between loss and profit.

Efficient methods for solving linear programming problems became available in the late 1930s. In 1939, Kantorovich presented a number of solutions to some problems related to production and transportation planning. During World War II, Koopmans contributed significantly to the solution of transportation problems. Kantorovich and Koopmans were awarded a Nobel Prize in Economics in 1975 for their work on the theory of optimal allocation of resources. In 1947, Dantzig developed a new method for solving linear programs, known today as the *simplex method* (see [34] for Dantzig's own treatment of the algorithm). In the following chapters we discuss the simplex method in detail. The simplex method is efficient and elegant and has been declared one of the 10 algorithms with the greatest influence on the development and practice of science and engineering in the twentieth century [40].

The simplex method has the undesirable property that in the worst case, the number of steps (and hence total time) required to find a solution grows exponentially with the number of variables. Thus, the simplex method is said to have *exponential worst-case complexity*. This led to an interest in devising algorithms for solving linear programs that have polynomial complexity—algorithms that find a solution in an amount of time that is bounded by a polynomial in the number of variables. Khachiyan, in 1979, was the first to devise such an algorithm. However, his algorithm gained more theoretical than practical interest. Then, in 1984, Karmarkar proposed a new linear program-

ming algorithm that has polynomial complexity and appears to solve some complicated real-world problems of scheduling, routing, and planning more efficiently than the simplex method. Karmarkar's work led to the development of many other nonsimplex methods commonly referred to as *interior-point methods*. This approach is currently still an active research area. For more details on Karmarkar's and related algorithms, see [42], [55], [71], [119], and [124]. Some basic ideas illustrating Khachiyan's and Karmarkar's algorithms are presented in Chapter 18.

## 15.2 Simple Examples of Linear Programs

Formally, a linear program is an optimization problem of the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . The vector inequality  $\mathbf{x} \geq \mathbf{0}$  means that each component of  $\mathbf{x}$  is nonnegative. Several variations of this problem are possible; for example, instead of minimizing, we can maximize, or the constraints may be in the form of inequalities, such as  $\mathbf{Ax} \geq \mathbf{b}$  or  $\mathbf{Ax} \leq \mathbf{b}$ . We also refer to these variations as linear programs. In fact, as we shall see later, these variations can all be rewritten into the standard form shown above.

The purpose of this section is to give some simple examples of linear programming problems illustrating the importance and the various applications of linear programming methods.

**Example 15.1** This example is adapted from [123]. A manufacturer produces four different products:  $X_1$ ,  $X_2$ ,  $X_3$ , and  $X_4$ . There are three inputs to this production process: labor in person-weeks, kilograms of raw material A, and boxes of raw material B. Each product has different input requirements. In determining each week's production schedule, the manufacturer cannot use more than the available amounts of labor and the two raw materials. The relevant information is presented in Table 15.1. Every production decision must satisfy the restrictions on the availability of inputs. These constraints can be written using the data in Table 15.1. In particular, we have

$$\begin{aligned} x_1 + 2x_2 + x_3 + 2x_4 &\leq 20 \\ 6x_1 + 5x_2 + 3x_3 + 2x_4 &\leq 100 \\ 3x_1 + 4x_2 + 9x_3 + 12x_4 &\leq 75. \end{aligned}$$

Because negative production levels are not meaningful, we must impose the following nonnegativity constraints on the production levels:

$$x_i \geq 0, \quad i = 1, 2, 3, 4.$$

**Table 15.1** Data for Example 15.1

Inputs	Product				Input Availabilities
	$X_1$	$X_2$	$X_3$	$X_4$	
Person-weeks	1	2	1	2	20
Kilograms of material A	6	5	3	2	100
Boxes of material B	3	4	9	12	75
Production levels	$x_1$	$x_2$	$x_3$	$x_4$	

Now, suppose that one unit of product  $X_1$  sells for \$6, and  $X_2$ ,  $X_3$ , and  $X_4$  sell for \$4, \$7, and \$5, respectively. Then, the total revenue for any production decision  $(x_1, x_2, x_3, x_4)$  is

$$f(x_1, x_2, x_3, x_4) = 6x_1 + 4x_2 + 7x_3 + 5x_4.$$

The problem is then to maximize  $f$  subject to the given constraints (the three inequalities and four nonnegativity constraints). Using vector notation with

$$\mathbf{x} = [x_1, x_2, x_3, x_4]^\top,$$

the problem can be written in the compact form

$$\begin{aligned} & \text{maximize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{c}^\top &= [6, 4, 7, 5], \\ \mathbf{A} &= \begin{bmatrix} 1 & 2 & 1 & 2 \\ 6 & 5 & 3 & 2 \\ 3 & 4 & 9 & 12 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 20 \\ 100 \\ 75 \end{bmatrix}. \end{aligned}$$

■

Another example that illustrates linear programming involves determining the most economical diet that satisfies the basic minimum requirements for good health.

**Example 15.2 Diet Problem.** This example is adapted from [88]. Assume that  $n$  different food types are available. The  $j$ th food sells at a price  $c_j$  per unit. In addition, there are  $m$  basic nutrients. To achieve a balanced diet, you must receive at least  $b_i$  units of the  $i$ th nutrient per day. Assume that each

unit of food  $j$  contains  $a_{ij}$  units of the  $i$ th nutrient. Denote by  $x_j$  the number of units of food  $j$  in the diet. The objective is to select the  $x_j$  to minimize the total cost of the diet:

$$\text{minimize } c_1x_1 + c_2x_2 + \cdots + c_nx_n$$

subject to the nutritional constraints

$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \geq b_1$$

$$a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \geq b_2$$

 $\vdots$ 

$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \geq b_m,$$

and the nonnegativity constraints

$$x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0.$$

In the more compact vector notation, this problem becomes

$$\text{minimize } \mathbf{c}^\top \mathbf{x}$$

$$\text{subject to } \mathbf{A}\mathbf{x} \geq \mathbf{b}$$

$$\mathbf{x} \geq \mathbf{0},$$

where  $\mathbf{x} = [x_1, x_2, \dots, x_n]^\top$  is an  $n$ -dimensional column vector,  $\mathbf{c}^\top$  is an  $n$ -dimensional row vector,  $\mathbf{A}$  is an  $m \times n$  matrix, and  $\mathbf{b}$  is an  $m$ -dimensional column vector. We call this problem the *diet problem* and will return to it in Chapter 17. ■

In the next example we consider a linear programming problem that arises in manufacturing.

**Example 15.3** A manufacturer produces two different products,  $X_1$  and  $X_2$ , using three machines:  $M_1$ ,  $M_2$ , and  $M_3$ . Each machine can be used for only a limited amount of time. Production times of each product on each machine are given in Table 15.2. The objective is to maximize the combined time of utilization of all three machines.

Every production decision must satisfy the constraints on the available time. These restrictions can be written down using data from Table 15.2. In particular, we have

$$x_1 + x_2 \leq 8,$$

$$x_1 + 3x_2 \leq 18,$$

$$2x_1 + x_2 \leq 14,$$

where  $x_1$  and  $x_2$  denote the production levels. The combined production time of all three machines is

$$f(x_1, x_2) = 4x_1 + 5x_2.$$

**Table 15.2** Data for Example 15.3

Machine	Production time (hours/unit)		(hours)
	$X_1$	$X_2$	
$M_1$	1	1	8
$M_2$	1	3	18
$M_3$	2	1	14
Total	4	5	

Thus, writing  $\mathbf{x} = [x_1, x_2]^\top$ , the problem in compact notation has the form

$$\begin{aligned} & \text{maximize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{c}^\top &= [4, 5], \\ \mathbf{A} &= \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 8 \\ 18 \\ 14 \end{bmatrix}. \end{aligned}$$

■

In the following example we discuss an application of linear programming in transportation.

**Example 15.4** A manufacturing company has plants in cities A, B, and C. The company produces and distributes its product to dealers in various cities. On a particular day, the company has 30 units of its product in A, 40 in B, and 30 in C. The company plans to ship 20 units to D, 20 to E, 25 to F, and 35 to G, following orders received from dealers. The transportation costs per unit of each product between the cities are given in Table 15.3. In the table, the quantities supplied and demanded appear at the right and along the bottom of the table. The quantities to be transported from the plants to different destinations are represented by the decision variables.

**Table 15.3** Data for Example 15.4

To From	D	E	F	G	Supply
A	\$7	\$10	\$14	\$8	30
B	\$7	\$11	\$12	\$6	40
C	\$5	\$8	\$15	\$9	30
Demand	20	20	25	35	100

This problem can be stated in the form

$$\begin{aligned} \text{minimize} \quad & 7x_{11} + 10x_{12} + 14x_{13} + 8x_{14} + 7x_{21} + 11x_{22} + 12x_{23} \\ & + 6x_{24} + 5x_{31} + 8x_{32} + 15x_{33} + 9x_{34} \end{aligned}$$

$$\text{subject to} \quad x_{11} + x_{12} + x_{13} + x_{14} = 30$$

$$x_{21} + x_{22} + x_{23} + x_{24} = 40$$

$$x_{31} + x_{32} + x_{33} + x_{34} = 30$$

$$x_{11} + x_{21} + x_{31} = 20$$

$$x_{12} + x_{22} + x_{32} = 20$$

$$x_{13} + x_{23} + x_{33} = 25$$

$$x_{14} + x_{24} + x_{34} = 35$$

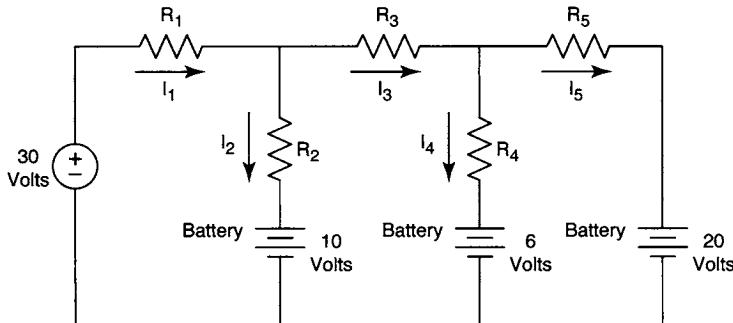
$$x_{11}, x_{12}, \dots, x_{34} \geq 0.$$

In this problem one of the constraint equations is redundant because it can be derived from the rest of the constraint equations. The mathematical formulation of the transportation problem is then in a linear programming form with twelve ( $3 \times 4$ ) decision variables and six ( $3 + 4 - 1$ ) linearly independent constraint equations. Obviously, we also require nonnegativity of the decision variables, since a negative shipment is impossible and does not have a valid interpretation. ■

Next, we give an example of a linear programming problem arising in electrical engineering.

**Example 15.5** This example is adapted from [100]. Figure 15.1 shows an electric circuit that is designed to use a 30-V source to charge 10-V, 6-V, and 20-V batteries connected in parallel. Physical constraints limit the currents  $I_1, I_2, I_3, I_4$ , and  $I_5$  to a maximum of 4 A, 3 A, 3 A, 2 A, and 2 A, respectively. In addition, the batteries must not be discharged; that is, the currents  $I_1, I_2, I_3, I_4$ , and  $I_5$  must not be negative. We wish to find the values of the currents  $I_1, \dots, I_5$  such that the total power transferred to the batteries is maximized.

The total power transferred to the batteries is the sum of the powers transferred to each battery and is given by  $10I_2 + 6I_4 + 20I_5$  W. From the circuit in



**Figure 15.1** Battery charger circuit for Example 15.5.

Figure 15.1, we observe that the currents satisfy the constraints  $I_1 = I_2 + I_3$  and  $I_3 = I_4 + I_5$ . Therefore, the problem can be posed as the following linear program:

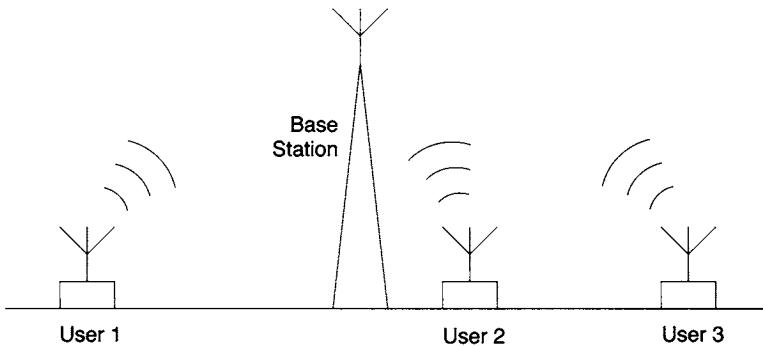
$$\begin{aligned} & \text{maximize} && 10I_2 + 6I_4 + 20I_5 \\ & \text{subject to} && I_1 = I_2 + I_3 \\ & && I_3 = I_4 + I_5 \\ & && I_1 \leq 4 \\ & && I_2 \leq 3 \\ & && I_3 \leq 3 \\ & && I_4 \leq 2 \\ & && I_5 \leq 2 \\ & && I_1, I_2, I_3, I_4, I_5 \geq 0. \end{aligned}$$

■

Finally, we present an example from wireless communications.

**Example 15.6** Consider the wireless communication system shown in Figure 15.2. There are  $n$  “mobile” users. For each  $i = 1, \dots, n$ , user  $i$  transmits a signal to the base station with power  $p_i$  and an attenuation factor of  $h_i$  (i.e., the actual signal power received at the base station from user  $i$  is  $h_i p_i$ ). When the base station is receiving from user  $i$ , the total power received from all other users is considered *interference* (i.e., the interference for user  $i$  is  $\sum_{j \neq i} h_j p_j$ ). For the communication with user  $i$  to be reliable, the signal-to-interference ratio must exceed a threshold  $\gamma_i$ , where the “signal” is the power received from user  $i$ .

We are interested in minimizing the total power transmitted by all users subject to having reliable communications for all users. We can formulate the



**Figure 15.2** Wireless communication system in Example 15.6.

problem as a linear programming problem of the form

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

We proceed as follows. The total power transmitted is  $p_1 + \dots + p_n$ . The signal-to-interference ratio for user  $i$  is

$$\frac{h_i p_i}{\sum_{j \neq i} h_j p_j}.$$

Hence, the problem can be written as

$$\begin{aligned} & \text{minimize } p_1 + \dots + p_n \\ & \text{subject to } \frac{h_i p_i}{\sum_{j \neq i} h_j p_j} \geq \gamma_i, \quad i = 1, \dots, n \\ & \quad p_1, \dots, p_n \geq 0. \end{aligned}$$

We can write the above as the linear programming problem

$$\begin{aligned} & \text{minimize } p_1 + \dots + p_n \\ & \text{subject to } h_i p_i - \gamma_i \sum_{j \neq i} h_j p_j \geq 0, \quad i = 1, \dots, n \\ & \quad p_1, \dots, p_n \geq 0. \end{aligned}$$

In matrix form, we have

$$\mathbf{c} = [1, \dots, 1]^\top$$

$$\mathbf{A} = \begin{bmatrix} h_1 & -\gamma_1 h_2 & \cdots & -\gamma_1 h_n \\ -\gamma_2 h_1 & h_2 & \cdots & -\gamma_2 h_n \\ \vdots & \ddots & \ddots & \vdots \\ -\gamma_n h_1 & -\gamma_n h_2 & \cdots & h_n \end{bmatrix}, \quad \mathbf{b} = \mathbf{0}.$$

■

For more examples of linear programming and their applications in a variety of engineering problems, we refer the reader to [1], [34], [35], [46], and [109]. For applications of linear programming to the design of control systems, see [33]. Linear programming also provides the basis for theoretical applications, as, for example, in matrix game theory (discussed in [18]).

### 15.3 Two-Dimensional Linear Programs

Many fundamental concepts of linear programming are easily illustrated in two-dimensional space. Therefore, we consider linear problems in  $\mathbb{R}^2$  before discussing general linear programming problems.

Consider the following linear program (adapted from [123]):

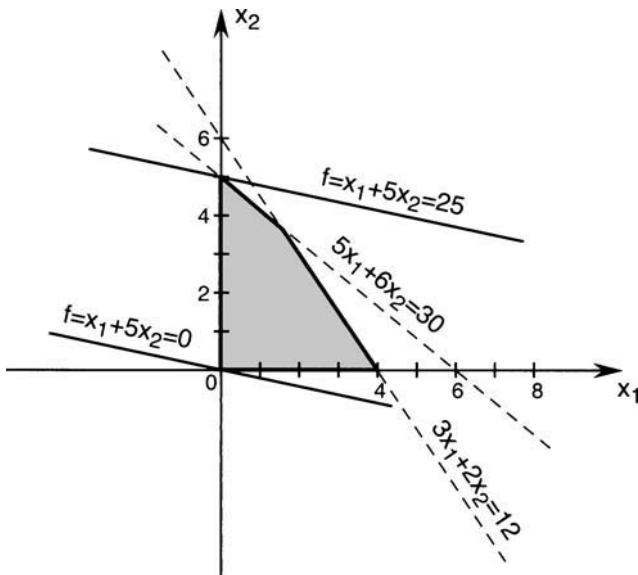
$$\begin{aligned} & \text{maximize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{x} = [x_1, x_2]^\top$  and

$$\mathbf{c}^\top = [1, 5],$$

$$\mathbf{A} = \begin{bmatrix} 5 & 6 \\ 3 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 30 \\ 12 \end{bmatrix}.$$

First, we note that the set of equations  $\{\mathbf{c}^\top \mathbf{x} = x_1 + 5x_2 = f, f \in \mathbb{R}\}$  specifies a family of straight lines in  $\mathbb{R}^2$ . Each member of this family can be obtained by setting  $f$  equal to some real number. Thus, for example,  $x_1 + 5x_2 = -5$ ,  $x_1 + 5x_2 = 0$ , and  $x_1 + 5x_2 = 3$  are three parallel lines belonging to the family. Now, suppose that we try to choose several values for  $x_1$  and  $x_2$  and observe how large we can make  $f$  while still satisfying the constraints on  $x_1$  and  $x_2$ . We first try  $x_1 = 1$  and  $x_2 = 3$ . This point satisfies the constraints. For this point,  $f = 16$ . If we now select  $x_1 = 0$  and  $x_2 = 5$ , then  $f = 25$  and this point yields a larger value for  $f$  than does  $\mathbf{x} = [1, 3]^\top$ . There are infinitely many points  $[x_1, x_2]^\top$  satisfying the constraints. Therefore, we need a better



**Figure 15.3** Geometric solution of a linear program in  $\mathbb{R}^2$ .

method than trial and error to solve the problem. In the following sections we develop a systematic approach that simplifies considerably the process of solving linear programming problems.

For the example above we can easily solve the problem using geometric arguments. First let us sketch the constraints in  $\mathbb{R}^2$ . The region of feasible points (the set of points  $\mathbf{x}$  satisfying the constraints  $A\mathbf{x} \leq \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ ) is depicted by the shaded region in Figure 15.3.

Geometrically, maximizing  $\mathbf{c}^\top \mathbf{x} = x_1 + 5x_2$  subject to the constraints can be thought of as finding the straight line  $f = x_1 + 5x_2$  that intersects the shaded region and has the largest  $f$ . The coordinates of the point of intersection will then yield a maximum value of  $\mathbf{c}^\top \mathbf{x}$ . In our example, the point  $[0, 5]^\top$  is the solution (see Figure 15.3).

**Example 15.7** Suppose that you are given two different types of concrete. The first type contains 30% cement, 40% gravel, and 30% sand (all percentages of weight). The second type contains 10% cement, 20% gravel, and 70% sand. The first type of concrete costs \$5 per pound and the second type costs \$1 per pound. How many pounds of each type of concrete should you buy and mix together so that your cost is minimized but you get a concrete mixture that has at least a total of 5 pounds of cement, 3 pounds of gravel, and 4 pounds of sand?

The problem can be represented as

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{c}^\top &= [5, 1], \\ \mathbf{A} &= \begin{bmatrix} 0.3 & 0.1 \\ 0.4 & 0.2 \\ 0.3 & 0.7 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 5 \\ 3 \\ 4 \end{bmatrix}. \end{aligned}$$

Using the graphical method described above, we get a solution of  $[0, 50]^\top$ , which means that we should purchase 50 pounds of the second type of concrete. (For a variation of this problem solved using a different method, see Example 12.1.) ■

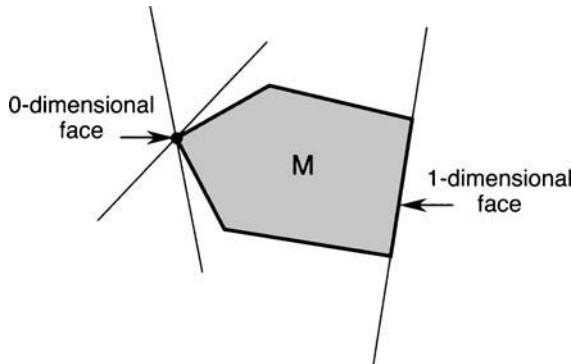
In some cases, when using the graphical method, there may be more than one point of intersection of the optimal straight line  $f = \mathbf{c}^\top \mathbf{x}$  with the boundary of the feasible region. In this case all of the intersection points will yield the same value for the objective function  $\mathbf{c}^\top \mathbf{x}$ , and therefore any one of them is a solution.

## 15.4 Convex Polyhedra and Linear Programming

The goal of linear programming is to minimize (or maximize) a linear objective function

$$\mathbf{c}^\top \mathbf{x} = c_1 x_1 + c_2 x_2 + \cdots + c_n x_n$$

subject to constraints that are represented by linear equalities and/or inequalities. For the time being, let us consider only constraints of the form  $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ . In this section we discuss linear programs from a geometric point of view (for a review of geometric concepts used in the section, see Chapter 4). The set of points satisfying these constraints can be represented as the intersection of a finite number of closed half-spaces. Thus, the constraints define a convex polytope. We assume, for simplicity, that this polytope is nonempty and bounded. In other words, the equations of constraints define a polyhedron  $M$  in  $\mathbb{R}^n$ . Let  $H$  be a hyperplane of support of this polyhedron. If the dimension of  $M$  is less than  $n$ , then the set of all points common to the hyperplane  $H$  and the polyhedron  $M$  coincides with  $M$ . If the dimension of  $M$  is equal to  $n$ , then the set of all points common to the hyperplane  $H$  and the polyhedron  $M$  is a face of the polyhedron. If this face is  $(n - 1)$ -dimensional, then there exists only one hyperplane of support, namely, the carrier of this face. If the



**Figure 15.4** Hyperplanes of support at different boundary points of the polyhedron  $M$ .

dimension of the face is less than  $n - 1$ , then there exist an infinite number of hyperplanes of support whose intersection with this polyhedron yields this face (see Figure 15.4).

The goal of our linear programming problem is to maximize a linear objective function  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x} = c_1 x_1 + \cdots + c_n x_n$  on the convex polyhedron  $M$ . Next, let  $H$  be the hyperplane defined by the equation

$$\mathbf{c}^\top \mathbf{x} = 0.$$

Draw a hyperplane of support  $\tilde{H}$  to the polyhedron  $M$ , which is parallel to  $H$  and positioned such that the vector  $\mathbf{c}$  points in the direction of the half-space that does not contain  $M$  (see Figure 15.5). The equation of the hyperplane  $\tilde{H}$  has the form

$$\mathbf{c}^\top \mathbf{x} = \beta,$$

and for all  $\mathbf{x} \in M$  we have  $\mathbf{c}^\top \mathbf{x} \leq \beta$ . Denote by  $\tilde{M}$  the convex polyhedron that is the intersection of the hyperplane of support  $\tilde{H}$  with the polyhedron  $M$ . We now show that  $f$  is constant on  $\tilde{M}$  and that  $\tilde{M}$  is the set of all points in  $M$  for which  $f$  attains its maximum value. To this end, let  $\mathbf{y}$  and  $\mathbf{z}$  be two arbitrary points in  $\tilde{M}$ . This implies that both  $\mathbf{y}$  and  $\mathbf{z}$  belong to  $\tilde{H}$ . Hence,

$$f(\mathbf{y}) = \mathbf{c}^\top \mathbf{y} = \beta = \mathbf{c}^\top \mathbf{z} = f(\mathbf{z}),$$

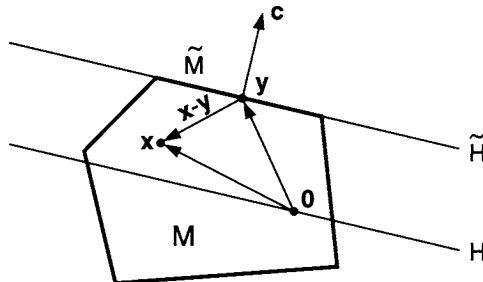
which means that  $f$  is constant on  $\tilde{M}$ .

Let  $\mathbf{y}$  be a point of  $\tilde{M}$ , and let  $\mathbf{x}$  be a point of  $M \setminus \tilde{M}$ ; that is,  $\mathbf{x}$  is a point of  $M$  that does not belong to  $\tilde{M}$  (see Figure 15.5). Then,

$$\mathbf{c}^\top \mathbf{x} < \beta = \mathbf{c}^\top \mathbf{y},$$

which implies that

$$f(\mathbf{x}) < f(\mathbf{y}).$$



**Figure 15.5** Maximization of a linear function on the polyhedron  $M$ .

Thus, the values of  $f$  at the points of  $M$  that do not belong to  $\tilde{M}$  are smaller than the values at points of  $\tilde{M}$ . Hence,  $f$  achieves its maximum on  $M$  at points in  $\tilde{M}$ .

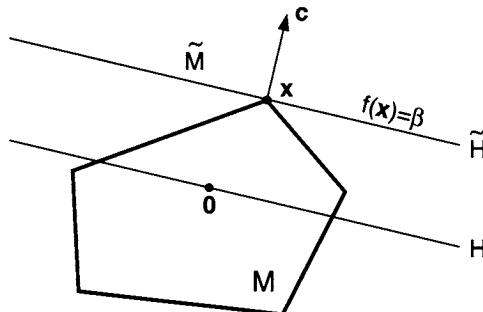
It may happen that  $\tilde{M}$  contains only a single point, in which case  $f$  achieves its maximum at a unique point. This occurs when the the hyperplane of support passes through an extreme point of  $M$  (see Figure 15.6).

### 15.5 Standard Form Linear Programs

We refer to a linear program of the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \end{aligned}$$

as a linear program in *standard form*. Here  $\mathbf{A}$  is an  $m \times n$  matrix composed of real entries,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ . Without loss of generality, we assume that



**Figure 15.6** Unique maximum point of  $f$  on the polyhedron  $M$ .

$\mathbf{b} \geq \mathbf{0}$ . If a component of  $\mathbf{b}$  is negative, say the  $i$ th component, we multiply the  $i$ th constraint by  $-1$  to obtain a positive right-hand side.

Theorems and solution techniques for linear programs are usually stated for problems in standard form. Other forms of linear programs can be converted to the standard form, as we now show. If a linear program is in the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

then by introducing *surplus variables*  $y_i$ , we can convert the original problem into the standard form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{in}x_n - y_i = b_i, \quad i = 1, \dots, m \\ & && x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0 \\ & && y_1 \geq 0, y_2 \geq 0, \dots, y_m \geq 0. \end{aligned}$$

In more compact notation, the formulation above can be represented as

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{Ax} - \mathbf{I}_m \mathbf{y} = [\mathbf{A}, -\mathbf{I}_m] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \mathbf{y} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{I}_m$  is the  $m \times m$  identity matrix.

If, on the other hand, the constraints have the form

$$\begin{aligned} & \mathbf{Ax} \leq \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

then we introduce *slack variables*  $y_i$  to convert the constraints into the form

$$\begin{aligned} & \mathbf{Ax} + \mathbf{I}_m \mathbf{y} = [\mathbf{A}, \mathbf{I}_m] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}, \mathbf{y} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{y}$  is the vector of slack variables. Note that neither surplus nor slack variables contribute to the objective function  $\mathbf{c}^\top \mathbf{x}$ .

At first glance, it may appear that the two problems

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{Ax} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \end{aligned}$$

and

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{Ax} - \mathbf{I}_m \mathbf{y} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \\ & && \mathbf{y} \geq \mathbf{0} \end{aligned}$$

are different, in that the first problem refers to the intersection of half-spaces in the  $n$ -dimensional space, whereas the second problem refers to an intersection of half-spaces and hyperplanes in the  $(n+m)$ -dimensional space. It turns out that both formulations are algebraically equivalent in the sense that a solution to one of the problems implies a solution to the other. To illustrate this equivalence, we consider the following examples.

**Example 15.8** Suppose that we are given the inequality constraint

$$x_1 \leq 7.$$

We convert this to an equality constraint by introducing a slack variable  $x_2 \geq 0$  to obtain

$$\begin{aligned} x_1 + x_2 &= 7 \\ x_2 &\geq 0. \end{aligned}$$

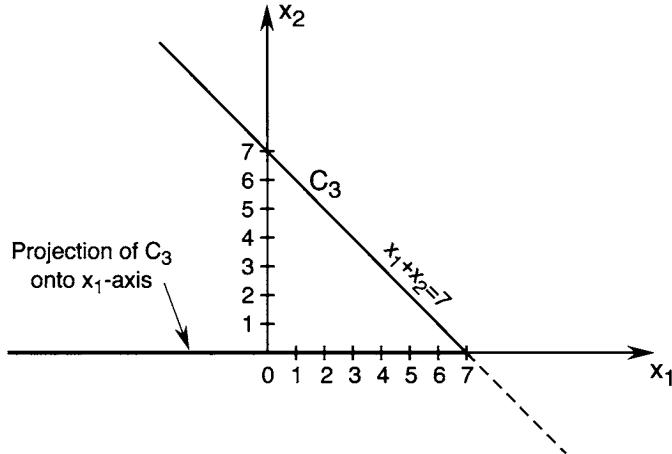
Consider the sets  $C_1 = \{x_1 : x_1 \leq 7\}$  and  $C_2 = \{x_1 : x_1 + x_2 = 7, x_2 \geq 0\}$ . Are the sets  $C_1$  and  $C_2$  equal? It is clear that indeed they are; in this example, we give a geometric interpretation for their equality. Consider a third set  $C_3 = \{[x_1, x_2]^\top : x_1 + x_2 = 7, x_2 \geq 0\}$ . From Figure 15.7 we can see that the set  $C_3$  consists of all points on the line to the left and above the point of intersection of the line with the  $x_1$ -axis. This set, being a subset of  $\mathbb{R}^2$ , is of course not the same set as the set  $C_1$  (a subset of  $\mathbb{R}$ ). However, we can project the set  $C_3$  onto the  $x_1$ -axis (see Figure 15.7). We can associate with each point  $x_1 \in C_1$  a point  $[x_1, 0]^\top$  on the orthogonal projection of  $C_3$  onto the  $x_1$ -axis, and vice versa. Note that  $C_2 = \{x_1 : [x_1, x_2]^\top \in C_3\} = C_1$ . ■

**Example 15.9** Consider the inequality constraints

$$\begin{aligned} a_1 x_1 + a_2 x_2 &\leq b \\ x_1, x_2 &\geq 0, \end{aligned}$$

where  $a_1$ ,  $a_2$ , and  $b$  are positive numbers. Again, we introduce a slack variable  $x_3 \geq 0$  to get

$$\begin{aligned} a_1 x_1 + a_2 x_2 + x_3 &= b \\ x_1, x_2, x_3 &\geq 0. \end{aligned}$$



**Figure 15.7** Projection of the set  $C_3$  onto the  $x_1$ -axis.

Define the sets

$$\begin{aligned}C_1 &= \{[x_1, x_2]^\top : a_1 x_1 + a_2 x_2 \leq b, x_1, x_2 \geq 0\}, \\C_2 &= \{[x_1, x_2]^\top : a_1 x_1 + a_2 x_2 + x_3 = b, x_1, x_2, x_3 \geq 0\}, \\C_3 &= \{[x_1, x_2, x_3]^\top : a_1 x_1 + a_2 x_2 + x_3 = b, x_1, x_2, x_3 \geq 0\}.\end{aligned}$$

We again see that  $C_3$  is not the same as  $C_1$ . However, the orthogonal projection of  $C_3$  onto the  $(x_1, x_2)$ -plane allows us to associate the resulting set with the set  $C_1$ . We associate the points  $[x_1, x_2, 0]^\top$  resulting from the orthogonal projection of  $C_3$  onto the  $(x_1, x_2)$ -plane with the points in  $C_1$  (see Figure 15.8). Note that  $C_2 = \{[x_1, x_2]^\top : [x_1, x_2, x_3]^\top \in C_3\} = C_1$ . ■

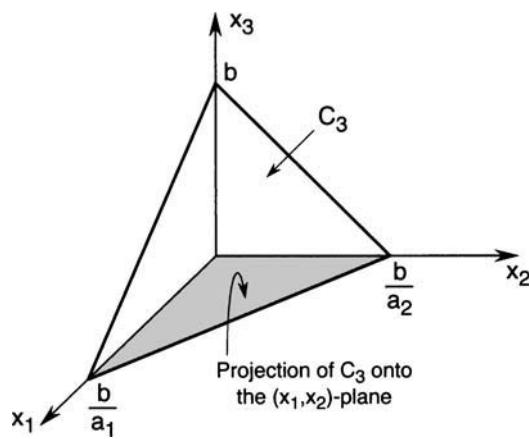
**Example 15.10** Suppose that we wish to maximize

$$f(x_1, x_2) = c_1 x_1 + c_2 x_2$$

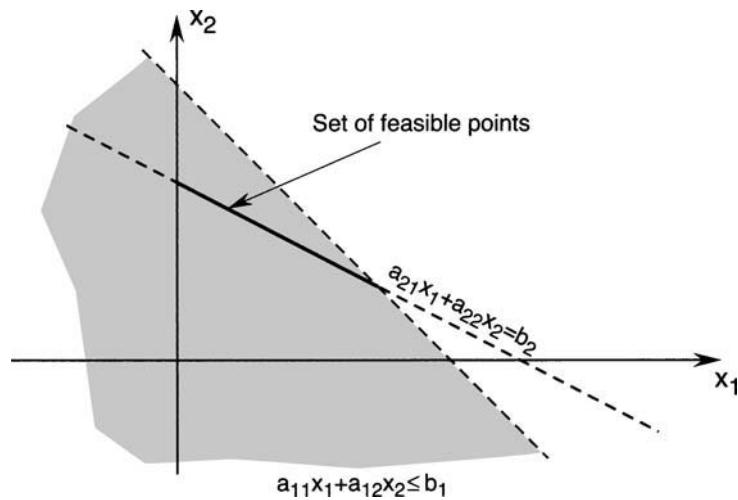
subject to the constraints

$$\begin{aligned}a_{11} x_1 + a_{12} x_2 &\leq b_1 \\a_{21} x_1 + a_{22} x_2 &= b_2 \\x_1, x_2 &\geq 0,\end{aligned}$$

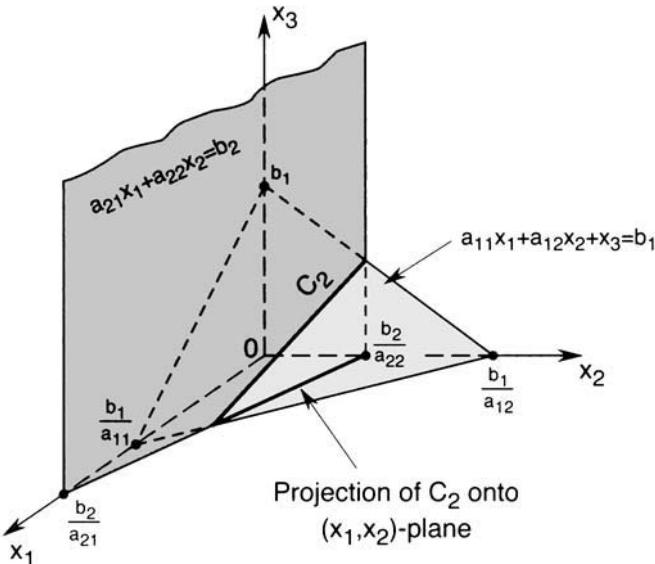
where, for simplicity, we assume that each  $a_{ij} > 0$  and  $b_1, b_2 \geq 0$ . The set of feasible points is depicted in Figure 15.9. Let  $C_1 \subset \mathbb{R}^2$  be the set of points satisfying the constraints.



**Figure 15.8** Projection of the set  $C_3$  onto the  $(x_1, x_2)$ -plane.



**Figure 15.9** The feasible set for Example 15.10.



**Figure 15.10** Projection of  $C_2$  onto the  $(x_1, x_2)$ -plane.

Introducing a slack variable, we convert the constraints into standard form:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + x_3 &= b_1 \\ a_{21}x_1 + a_{22}x_2 &= b_2 \\ x_i &\geq 0, \quad i = 1, 2, 3. \end{aligned}$$

Let  $C_2 \subset \mathbb{R}^3$  be the set of points satisfying the constraints. As illustrated in Figure 15.10, this set is a line segment (in  $\mathbb{R}^3$ ). We now project  $C_2$  onto the  $(x_1, x_2)$ -plane. The projected set consists of the points  $[x_1, x_2, 0]^\top$ , with  $[x_1, x_2, x_3]^\top \in C_2$  for some  $x_3 \geq 0$ . In Figure 15.10 this set is marked by a heavy line in the  $(x_1, x_2)$ -plane. We can associate the points on the projection with the corresponding points in the set  $C_1$ . ■

In the following example we convert an optimization problem into a standard form linear programming problem.

**Example 15.11** Consider the following optimization problem

$$\begin{aligned} &\text{maximize} && x_2 - x_1 \\ &\text{subject to} && 3x_1 = x_2 - 5 \\ & && |x_2| \leq 2 \\ & && x_1 \leq 0. \end{aligned}$$

To convert the problem into a standard form linear programming problem, we perform the following steps:

1. Change the objective function to: minimize  $x_1 - x_2$ .
2. Substitute  $x_1 = -x'_1$ .
3. Write  $|x_2| \leq 2$  as  $x_2 \leq 2$  and  $-x_2 \leq 2$ .
4. Introduce slack variables  $x_3$  and  $x_4$ , and convert the inequalities above to  $x_2 + x_3 = 2$  and  $-x_2 + x_4 = 2$ .
5. Write  $x_2 = u - v$ ,  $u, v \geq 0$ .

Hence, we obtain

$$\begin{aligned} & \text{minimize} && -x'_1 - u + v \\ & \text{subject to} && 3x'_1 + u - v = 5 \\ & && u - v + x_3 = 2 \\ & && v - u + x_4 = 2 \\ & && x'_1, u, v, x_3, x_4 \geq 0. \end{aligned}$$

■

## 15.6 Basic Solutions

We have seen in Section 15.5 that any linear programming problem involving inequalities can be converted to *standard form*, that is, a problem involving linear equations with nonnegative variables:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ , and  $\mathbf{b} \geq \mathbf{0}$ . In the following discussion we only consider linear programming problems in standard form.

Consider the system of equalities

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where  $\text{rank } \mathbf{A} = m$ . In dealing with this system of equations, we frequently need to consider a subset of columns of the matrix  $\mathbf{A}$ . For convenience, we often reorder the columns of  $\mathbf{A}$  so that the columns we are interested in appear first. Specifically, let  $\mathbf{B}$  be a square matrix whose columns are  $m$  linearly independent columns of  $\mathbf{A}$ . If necessary, we reorder the columns of  $\mathbf{A}$  so that the columns in  $\mathbf{B}$  appear first:  $\mathbf{A}$  has the form  $\mathbf{A} = [\mathbf{B}, \mathbf{D}]$ , where

$\mathbf{D}$  is an  $m \times (n - m)$  matrix whose columns are the remaining columns of  $\mathbf{A}$ . The matrix  $\mathbf{B}$  is nonsingular, and thus we can solve the equation

$$\mathbf{B}\mathbf{x}_B = \mathbf{b}$$

for the  $m$ -vector  $\mathbf{x}_B$ . The solution is  $\mathbf{x}_B = \mathbf{B}^{-1}\mathbf{b}$ . Let  $\mathbf{x}$  be the  $n$ -vector whose first  $m$  components are equal to  $\mathbf{x}_B$  and the remaining components are equal to zero; that is,  $\mathbf{x} = [\mathbf{x}_B^\top, \mathbf{0}^\top]^\top$ . Then,  $\mathbf{x}$  is a solution to  $\mathbf{Ax} = \mathbf{b}$ .

**Definition 15.1** We call  $[\mathbf{x}_B^\top, \mathbf{0}^\top]^\top$  a *basic solution* to  $\mathbf{Ax} = \mathbf{b}$  with respect to the basis  $\mathbf{B}$ . We refer to the components of the vector  $\mathbf{x}_B$  as *basic variables* and the columns of  $\mathbf{B}$  as *basic columns*.

If some of the basic variables of a basic solution are zero, then the basic solution is said to be a *degenerate basic solution*.

A vector  $\mathbf{x}$  satisfying  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ , is said to be a *feasible solution*.

A feasible solution that is also basic is called a *basic feasible solution*.

If the basic feasible solution is a degenerate basic solution, then it is called a *degenerate basic feasible solution*. ■

Note that in any basic feasible solution,  $\mathbf{x}_B \geq \mathbf{0}$ .

**Example 15.12** Consider the equation  $\mathbf{Ax} = \mathbf{b}$  with

$$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4] = \begin{bmatrix} 1 & 1 & -1 & 4 \\ 1 & -2 & -1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 8 \\ 2 \end{bmatrix},$$

where  $\mathbf{a}_i$  denotes the  $i$ th column of the matrix  $\mathbf{A}$ .

Then,  $\mathbf{x} = [6, 2, 0, 0]^\top$  is a basic feasible solution with respect to the basis  $\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_2]$ ,  $\mathbf{x} = [0, 0, 0, 2]^\top$  is a degenerate basic feasible solution with respect to the basis  $\mathbf{B} = [\mathbf{a}_3, \mathbf{a}_4]$  (as well as  $[\mathbf{a}_1, \mathbf{a}_4]$  and  $[\mathbf{a}_2, \mathbf{a}_4]$ ),  $\mathbf{x} = [3, 1, 0, 1]^\top$  is a feasible solution that is not basic, and  $\mathbf{x} = [0, 2, -6, 0]^\top$  is a basic solution with respect to the basis  $\mathbf{B} = [\mathbf{a}_2, \mathbf{a}_3]$ , but is not feasible. ■

**Example 15.13** As another example, consider the system of linear equations  $\mathbf{Ax} = \mathbf{b}$ , where

$$\mathbf{A} = \begin{bmatrix} 2 & 3 & -1 & -1 \\ 4 & 1 & 1 & -2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ 9 \end{bmatrix}.$$

We now find all solutions of this system. Note that every solution  $\mathbf{x}$  of  $\mathbf{Ax} = \mathbf{b}$  has the form  $\mathbf{x} = \mathbf{v} + \mathbf{h}$ , where  $\mathbf{v}$  is a particular solution of  $\mathbf{Ax} = \mathbf{b}$  and  $\mathbf{h}$  is a solution to  $\mathbf{Ax} = \mathbf{0}$ .

We form the augmented matrix  $[\mathbf{A}, \mathbf{b}]$  of the system:

$$[\mathbf{A}, \mathbf{b}] = \begin{bmatrix} 2 & 3 & -1 & -1 & -1 \\ 4 & 1 & 1 & -2 & 9 \end{bmatrix}.$$

Using elementary row operations, we transform this matrix into the form (see Chapter 16) given by

$$\begin{bmatrix} 1 & 0 & \frac{2}{5} & -\frac{1}{2} & \frac{14}{5} \\ 0 & 1 & -\frac{3}{5} & 0 & -\frac{11}{5} \end{bmatrix}.$$

The corresponding system of equations is given by

$$\begin{aligned} x_1 + \frac{2}{5}x_3 - \frac{1}{2}x_4 &= \frac{14}{5} \\ x_2 - \frac{3}{5}x_3 &= -\frac{11}{5}. \end{aligned}$$

Solving for the leading unknowns  $x_1$  and  $x_2$ , we obtain

$$\begin{aligned} x_1 &= \frac{14}{15} - \frac{2}{5}x_3 + \frac{1}{2}x_4 \\ x_2 &= -\frac{11}{5} + \frac{3}{5}x_3, \end{aligned}$$

where  $x_3$  and  $x_4$  are arbitrary real numbers. If  $[x_1, x_2, x_3, x_4]^\top$  is a solution, then we have

$$\begin{aligned} x_1 &= \frac{14}{5} - \frac{2}{5}s + \frac{1}{2}t, \\ x_2 &= -\frac{11}{5} + \frac{3}{5}s, \\ x_3 &= s, \\ x_4 &= t, \end{aligned}$$

where we have substituted  $s$  and  $t$  for  $x_3$  and  $x_4$ , respectively, to indicate that they are arbitrary real numbers.

Using vector notation, we may write the system of equations above as

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} \frac{14}{5} \\ -\frac{11}{5} \\ 0 \\ 0 \end{bmatrix} + s \begin{bmatrix} -\frac{2}{5} \\ \frac{3}{5} \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} \frac{1}{2} \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

Note that we have infinitely many solutions, parameterized by  $s, t \in \mathbb{R}$ . For the choice  $s = t = 0$  we obtain a particular solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , given by

$$\mathbf{v} = \begin{bmatrix} \frac{14}{5} \\ -\frac{11}{5} \\ 0 \\ 0 \end{bmatrix}.$$

Any other solution has the form  $\mathbf{v} + \mathbf{h}$ , where

$$\mathbf{h} = s \begin{bmatrix} -\frac{2}{5} \\ \frac{3}{5} \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} \frac{1}{2} \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

The total number of possible basic solutions is at most

$$\binom{n}{m} = \frac{n!}{m!(n-m)!} = \frac{4!}{2!(4-2)!} = 6.$$

To find basic solutions that are feasible, we check each of the basic solutions for feasibility.

Our first candidate for a basic feasible solution is obtained by setting  $x_3 = x_4 = 0$ , which corresponds to the basis  $\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_2]$ . Solving  $\mathbf{Bx}_B = \mathbf{b}$ , we obtain  $\mathbf{x}_B = [14/5, -11/5]^\top$ , and hence  $\mathbf{x} = [14/5, -11/5, 0, 0]^\top$  is a basic solution that is not feasible.

For our second candidate basic feasible solution, we set  $x_2 = x_4 = 0$ . We have the basis  $\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_3]$ . Solving  $\mathbf{Bx}_B = \mathbf{b}$  yields  $\mathbf{x}_B = [4/3, 11/3]^\top$ . Hence,  $\mathbf{x} = [4/3, 0, 11/3, 0]^\top$  is a basic feasible solution.

A third candidate basic feasible solution is obtained by setting  $x_2 = x_3 = 0$ . However, the matrix

$$\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_4] = \begin{bmatrix} 2 & -1 \\ 4 & -2 \end{bmatrix}$$

is singular. Therefore,  $\mathbf{B}$  cannot be a basis, and we do not have a basic solution corresponding to  $\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_4]$ .

We get our fourth candidate for a basic feasible solution by setting  $x_1 = x_4 = 0$ . We have a basis  $\mathbf{B} = [\mathbf{a}_2, \mathbf{a}_3]$ , resulting in  $\mathbf{x} = [0, 2, 7, 0]^\top$ , which is a basic feasible solution.

Our fifth candidate for a basic feasible solution corresponds to setting  $x_1 = x_3 = 0$ , with the basis  $\mathbf{B} = [\mathbf{a}_2, \mathbf{a}_4]$ . This results in  $\mathbf{x} = [0, -11/5, 0, -28/5]^\top$ , which is a basic solution that is not feasible.

Finally, the sixth candidate for a basic feasible solution is obtained by setting  $x_1 = x_2 = 0$ . This results in the basis  $\mathbf{B} = [\mathbf{a}_3, \mathbf{a}_4]$ , and  $\mathbf{x} = [0, 0, 11/3, -8/3]^\top$ , which is a basic solution but is not feasible. ■

## 15.7 Properties of Basic Solutions

In this section we discuss the importance of basic feasible solutions in solving linear programming (LP) problems. We first prove the fundamental theorem of LP, which states that when solving an LP problem, we need only consider

basic feasible solutions. This is because the optimal value (if it exists) is always achieved at a basic feasible solution. We need the following definitions.

**Definition 15.2** Any vector  $\mathbf{x}$  that yields the minimum value of the objective function  $\mathbf{c}^\top \mathbf{x}$  over the set of vectors satisfying the constraints  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ , is said to be an *optimal feasible solution*.

An optimal feasible solution that is basic is said to be an *optimal basic feasible solution*. ■

**Theorem 15.1 Fundamental Theorem of LP.** Consider a linear program in standard form.

1. If there exists a feasible solution, then there exists a basic feasible solution.
2. If there exists an optimal feasible solution, then there exists an optimal basic feasible solution. □

*Proof.* We first prove part 1. Suppose that  $\mathbf{x} = [x_1, \dots, x_n]^\top$  is a feasible solution and it has  $p$  positive components. Without loss of generality, we can assume that the first  $p$  components are positive, whereas the remaining components are zero. Then, in terms of the columns of  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_p, \dots, \mathbf{a}_n]$ , this solution satisfies

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_p\mathbf{a}_p = \mathbf{b}.$$

There are now two cases to consider.

*Case 1:* If  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$  are linearly independent, then  $p \leq m$ . If  $p = m$ , then the solution  $\mathbf{x}$  is basic and the proof is done. If, on the other hand,  $p < m$ , then, since  $\text{rank } \mathbf{A} = m$ , we can find  $m - p$  columns of  $\mathbf{A}$  from the remaining  $n - p$  columns so that the resulting set of  $m$  columns forms a basis. Hence, the solution  $\mathbf{x}$  is a (degenerate) basic feasible solution corresponding to the basis above.

*Case 2:* Assume that  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$  are linearly dependent. Then, there exist numbers  $y_i$ ,  $i = 1, \dots, p$ , not all zero, such that

$$y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \cdots + y_p\mathbf{a}_p = \mathbf{0}.$$

We can assume that there exists at least one  $y_i$  that is positive, for if all the  $y_i$  are nonpositive, we can multiply the equation above by  $-1$ . Multiply the equation by a scalar  $\varepsilon$  and subtract the resulting equation from  $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_p\mathbf{a}_p = \mathbf{b}$  to obtain

$$(x_1 - \varepsilon y_1)\mathbf{a}_1 + (x_2 - \varepsilon y_2)\mathbf{a}_2 + \cdots + (x_p - \varepsilon y_p)\mathbf{a}_p = \mathbf{b}.$$

Let

$$\mathbf{y} = [y_1, \dots, y_p, 0, \dots, 0]^\top.$$

Then, for any  $\varepsilon$  we can write

$$\mathbf{A}[\mathbf{x} - \varepsilon \mathbf{y}] = \mathbf{b}.$$

Let  $\varepsilon = \min\{x_i/y_i : i = 1, \dots, p, y_i > 0\}$ . Then, the first  $p$  components of  $\mathbf{x} - \varepsilon \mathbf{y}$  are nonnegative, and at least one of these components is zero. We then have a feasible solution with at most  $p - 1$  positive components. We can repeat this process until we get linearly independent columns of  $\mathbf{A}$ , after which we are back to case 1. Therefore, part 1 is proved.

We now prove part 2. Suppose that  $\mathbf{x} = [x_1, \dots, x_n]^\top$  is an optimal feasible solution and only the first  $p$  variables are nonzero. Then, we have two cases to consider. The first case is exactly the same as in part 1. The second case follows the same arguments as in part 1, but in addition we must show that  $\mathbf{x} - \varepsilon \mathbf{y}$  is optimal for any  $\varepsilon$ . We do this by showing that  $\mathbf{c}^\top \mathbf{y} = 0$ . To this end, assume that  $\mathbf{c}^\top \mathbf{y} \neq 0$ . Note that for  $\varepsilon$  of sufficiently small magnitude ( $|\varepsilon| \leq \min\{|x_i/y_i| : i = 1, \dots, p, y_i \neq 0\}$ ), the vector  $\mathbf{x} - \varepsilon \mathbf{y}$  is feasible. We can choose  $\varepsilon$  such that  $\mathbf{c}^\top \mathbf{x} > \mathbf{c}^\top \mathbf{x} - \varepsilon \mathbf{c}^\top \mathbf{y} = \mathbf{c}^\top (\mathbf{x} - \varepsilon \mathbf{y})$ . This contradicts the optimality of  $\mathbf{x}$ . We can now use the procedure from part 1 to obtain an optimal basic feasible solution from a given optimal feasible solution. ■

**Example 15.14** Consider the system of equations given in Example 15.13. Find a nonbasic feasible solution to this system and use the method in the proof of the fundamental theorem of LP to find a basic feasible solution.

Recall that solutions for the system given in Example 15.13 have the form

$$\mathbf{x} = \begin{bmatrix} \frac{14}{5} \\ -\frac{11}{5} \\ 0 \\ 0 \end{bmatrix} + s \begin{bmatrix} -\frac{2}{5} \\ \frac{3}{5} \\ 1 \\ 0 \end{bmatrix} + t \begin{bmatrix} \frac{1}{2} \\ 0 \\ 0 \\ 1 \end{bmatrix},$$

where  $s, t \in \mathbb{R}$ . Note that if  $s = 4$  and  $t = 0$ , then

$$\mathbf{x}_0 = \begin{bmatrix} \frac{6}{5} \\ \frac{1}{5} \\ \frac{4}{5} \\ 0 \end{bmatrix}$$

is a nonbasic feasible solution.

There are constants  $y_i$ ,  $i = 1, 2, 3$ , such that

$$y_1 \mathbf{a}_1 + y_2 \mathbf{a}_2 + y_3 \mathbf{a}_3 = \mathbf{0}.$$

For example, let

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

$$y_2 = \frac{3}{5},$$

$$y_3 = 1.$$

Note that

$$\mathbf{A}(\mathbf{x}_0 - \varepsilon \mathbf{y}) = \mathbf{b},$$

where

$$\mathbf{y} = \begin{bmatrix} -\frac{2}{5} \\ \frac{3}{5} \\ 1 \\ 0 \end{bmatrix}.$$

If  $\varepsilon = 1/3$ , then

$$\mathbf{x}_1 = \mathbf{x}_0 - \varepsilon \mathbf{y} = \begin{bmatrix} \frac{4}{3} \\ 0 \\ \frac{11}{3} \\ 0 \end{bmatrix}$$

is a basic feasible solution. ■

Observe that the fundamental theorem of LP reduces the task of solving a linear programming problem to that of searching over a finite number of basic feasible solutions. That is, we need only check basic feasible solutions for optimality. As mentioned before, the total number of basic solutions is at most

$$\binom{n}{m} = \frac{n!}{m!(n-m)!}.$$

Although this number is finite, it may be quite large. For example, if  $m = 5$  and  $n = 50$ , then

$$\binom{n}{m} = \binom{50}{5} = 2,118,760.$$

This is potentially the number of basic feasible solutions to be checked for optimality. Therefore, a more efficient method of solving linear programs is needed. To this end, in the next section we analyze a geometric interpretation of the fundamental theorem of LP. This leads us to the simplex method for solving linear programs, which we discuss in Chapter 16.

## 15.8 Geometric View of Linear Programs

Recall that a set  $\Theta \subset \mathbb{R}^n$  is said to be *convex* if, for every  $\mathbf{x}, \mathbf{y} \in \Theta$  and every real number  $\alpha$ ,  $0 < \alpha < 1$ , the point  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \Theta$ . In other words, a set is convex if given two points in the set, every point on the line segment joining these two points is also a member of the set.

Note that the set of points satisfying the constraints

$$\mathbf{Ax} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0}$$

is convex. To see this, let  $\mathbf{x}_1$  and  $\mathbf{x}_2$  satisfy the constraints, that is,  $A\mathbf{x}_i = \mathbf{b}$ ,  $\mathbf{x}_i \geq \mathbf{0}$ ,  $i = 1, 2$ . Then, for all  $\alpha \in (0, 1)$ ,  $A(\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) = \alpha A\mathbf{x}_1 + (1 - \alpha)A\mathbf{x}_2 = \mathbf{b}$ . Also, for  $\alpha \in (0, 1)$ , we have  $\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2 \geq \mathbf{0}$ .

Recall that a point  $\mathbf{x}$  in a convex set  $\Theta$  is said to be an *extreme point* of  $\Theta$  if there are no two distinct points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  in  $\Theta$  such that  $\mathbf{x} = \alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2$  for some  $\alpha \in (0, 1)$ . In other words, an extreme point is a point that does not lie strictly within the line segment connecting two other points of the set. Therefore, if  $\mathbf{x}$  is an extreme point, and  $\mathbf{x} = \alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2$  for some  $\mathbf{x}_1, \mathbf{x}_2 \in \Theta$  and  $\alpha \in (0, 1)$ , then  $\mathbf{x}_1 = \mathbf{x}_2$ . In the following theorem we show that extreme points of the constraint set are equivalent to basic feasible solutions.

**Theorem 15.2** *Let  $\Omega$  be the convex set consisting of all feasible solutions, that is, all  $n$ -vectors  $\mathbf{x}$  satisfying*

$$A\mathbf{x} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0},$$

*where  $A \in \mathbb{R}^{m \times n}$ ,  $m < n$ . Then,  $\mathbf{x}$  is an extreme point of  $\Omega$  if and only if  $\mathbf{x}$  is a basic feasible solution to  $A\mathbf{x} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ .*  $\square$

*Proof.*  $\Rightarrow$ : Suppose that  $\mathbf{x}$  satisfies  $A\mathbf{x} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ , and has  $p$  positive components. As before, without loss of generality, we can assume that the first  $p$  components are positive and the remaining components are zero. We have

$$x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_p\mathbf{a}_p = \mathbf{b}.$$

Let  $y_i$ ,  $i = 1, \dots, p$ , be numbers such that

$$y_1\mathbf{a}_1 + y_2\mathbf{a}_2 + \cdots + y_p\mathbf{a}_p = \mathbf{0}.$$

We show that each  $y_i = 0$ . To begin, multiply this equation by  $\varepsilon > 0$ , then add and subtract the result from the equation  $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_p\mathbf{a}_p = \mathbf{b}$  to get

$$\begin{aligned} (x_1 + \varepsilon y_1)\mathbf{a}_1 + (x_2 + \varepsilon y_2)\mathbf{a}_2 + \cdots + (x_p + \varepsilon y_p)\mathbf{a}_p &= \mathbf{b}, \\ (x_1 - \varepsilon y_1)\mathbf{a}_1 + (x_2 - \varepsilon y_2)\mathbf{a}_2 + \cdots + (x_p - \varepsilon y_p)\mathbf{a}_p &= \mathbf{b}. \end{aligned}$$

Because each  $x_i > 0$ ,  $\varepsilon > 0$  can be chosen such that each  $x_i + \varepsilon y_i, x_i - \varepsilon y_i \geq 0$  (e.g.,  $\varepsilon = \min\{|x_i/y_i| : i = 1, \dots, p, y_i \neq 0\}$ ). For such a choice of  $\varepsilon$ , the vectors

$$\begin{aligned} \mathbf{z}_1 &= [x_1 + \varepsilon y_1, x_2 + \varepsilon y_2, \dots, x_p + \varepsilon y_p, 0, \dots, 0]^\top, \\ \mathbf{z}_2 &= [x_1 - \varepsilon y_1, x_2 - \varepsilon y_2, \dots, x_p - \varepsilon y_p, 0, \dots, 0]^\top \end{aligned}$$

belong to  $\Omega$ . Observe that  $\mathbf{x} = \frac{1}{2}\mathbf{z}_1 + \frac{1}{2}\mathbf{z}_2$ . Because  $\mathbf{x}$  is an extreme point,  $\mathbf{z}_1 = \mathbf{z}_2$ . Hence, each  $y_i = 0$ , which implies that the  $\mathbf{a}_i$  are linearly independent.

$\Leftarrow$ : Let  $\mathbf{x} \in \Omega$  be a basic feasible solution. Let  $\mathbf{y}, \mathbf{z} \in \Omega$  be such that

$$\mathbf{x} = \alpha\mathbf{y} + (1 - \alpha)\mathbf{z}$$

for some  $\alpha \in (0, 1)$ . We show that  $\mathbf{y} = \mathbf{z}$  and conclude that  $\mathbf{x}$  is an extreme point. Because  $\mathbf{y}, \mathbf{z} \geq \mathbf{0}$ , and the last  $n - m$  components of  $\mathbf{x}$  are zero, the last  $n - m$  components of  $\mathbf{y}$  and  $\mathbf{z}$  are zero as well. Furthermore, since  $\mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{z} = \mathbf{b}$ ,

$$y_1\mathbf{a}_1 + \cdots + y_m\mathbf{a}_m = \mathbf{b}$$

and

$$z_1\mathbf{a}_1 + \cdots + z_m\mathbf{a}_m = \mathbf{b}.$$

Combining these two equations yields

$$(y_1 - z_1)\mathbf{a}_1 + \cdots + (y_m - z_m)\mathbf{a}_m = \mathbf{0}.$$

Because the columns  $\mathbf{a}_1, \dots, \mathbf{a}_m$  are linearly independent, we have  $y_i = z_i$ ,  $i = 1, \dots, m$ . Therefore,  $\mathbf{y} = \mathbf{z}$ , and hence  $\mathbf{x}$  is an extreme point of  $\Omega$ . ■

From Theorem 15.2 it follows that the set of extreme points of the constraint set  $\Omega = \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  is equal to the set of basic feasible solutions to  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ . Combining this observation with the fundamental theorem of LP (Theorem 15.1), we can see that in solving linear programming problems we need only examine the extreme points of the constraint set.

**Example 15.15** Consider the following LP problem:

$$\begin{aligned} & \text{maximize} && 3x_1 + 5x_2 \\ & \text{subject to} && x_1 + 5x_2 \leq 40 \\ & && 2x_1 + x_2 \leq 20 \\ & && x_1 + x_2 \leq 12 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

We introduce slack variables  $x_3, x_4, x_5$  to convert this LP problem into standard form:

$$\begin{aligned} & \text{minimize} && -3x_1 - 5x_2 \\ & \text{subject to} && x_1 + 5x_2 + x_3 = 40 \\ & && 2x_1 + x_2 + x_4 = 20 \\ & && x_1 + x_2 + x_5 = 12 \\ & && x_1, \dots, x_5 \geq 0. \end{aligned}$$

In the remainder of the example we consider only the problem in standard form. We can represent the constraints above as

$$\begin{aligned} x_1 \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} + x_2 \begin{bmatrix} 5 \\ 1 \\ 1 \end{bmatrix} + x_3 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + x_4 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + x_5 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} &= \begin{bmatrix} 40 \\ 20 \\ 12 \end{bmatrix}, \\ x_1, \dots, x_5 \geq 0, \end{aligned}$$

that is,  $x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + x_3\mathbf{a}_3 + x_4\mathbf{a}_4 + x_5\mathbf{a}_5 = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ . Note that

$$\mathbf{x} = [0, 0, 40, 20, 12]^\top$$

is a feasible solution. But for this  $\mathbf{x}$ , the value of the objective function is zero. We already know that the minimum of the objective function (if it exists) is achieved at an extreme point of the constraint set  $\Omega$  defined by the constraints. The point  $[0, 0, 40, 20, 12]^\top$  is an extreme point of the set of feasible solutions, but it turns out that it does not minimize the objective function. Therefore, we need to seek the solution among the other extreme points. To do this we move from one extreme point to an adjacent extreme point such that the value of the objective function decreases. Here, we define two extreme points to be adjacent if the corresponding basic columns differ by only one vector. We begin with  $\mathbf{x} = [0, 0, 40, 20, 12]^\top$ . We have

$$0\mathbf{a}_1 + 0\mathbf{a}_2 + 40\mathbf{a}_3 + 20\mathbf{a}_4 + 12\mathbf{a}_5 = \mathbf{b}.$$

To select an adjacent extreme point, let us choose to include  $\mathbf{a}_1$  as a basic column in the new basis. We need to remove either  $\mathbf{a}_3$ ,  $\mathbf{a}_4$ , or  $\mathbf{a}_5$  from the old basis. We proceed as follows. We first express  $\mathbf{a}_1$  as a linear combination of the old basic columns:

$$\mathbf{a}_1 = 1\mathbf{a}_3 + 2\mathbf{a}_4 + 1\mathbf{a}_5.$$

Multiplying both sides of this equation by  $\varepsilon_1 > 0$ , we get

$$\varepsilon_1\mathbf{a}_1 = \varepsilon_1\mathbf{a}_3 + 2\varepsilon_1\mathbf{a}_4 + \varepsilon_1\mathbf{a}_5.$$

We now add this equation to the equation  $0\mathbf{a}_1 + 0\mathbf{a}_2 + 40\mathbf{a}_3 + 20\mathbf{a}_4 + 12\mathbf{a}_5 = \mathbf{b}$ . Collecting terms yields

$$\varepsilon_1\mathbf{a}_1 + 0\mathbf{a}_2 + (40 - \varepsilon_1)\mathbf{a}_3 + (20 - 2\varepsilon_1)\mathbf{a}_4 + (12 - \varepsilon_1)\mathbf{a}_5 = \mathbf{b}.$$

We want to choose  $\varepsilon_1$  in such a way that each of the coefficients above is nonnegative and at the same time, one of the coefficients  $\mathbf{a}_3$ ,  $\mathbf{a}_4$ , or  $\mathbf{a}_5$  becomes zero. Clearly,  $\varepsilon_1 = 10$  does the job. The result is

$$10\mathbf{a}_1 + 30\mathbf{a}_3 + 2\mathbf{a}_5 = \mathbf{b}.$$

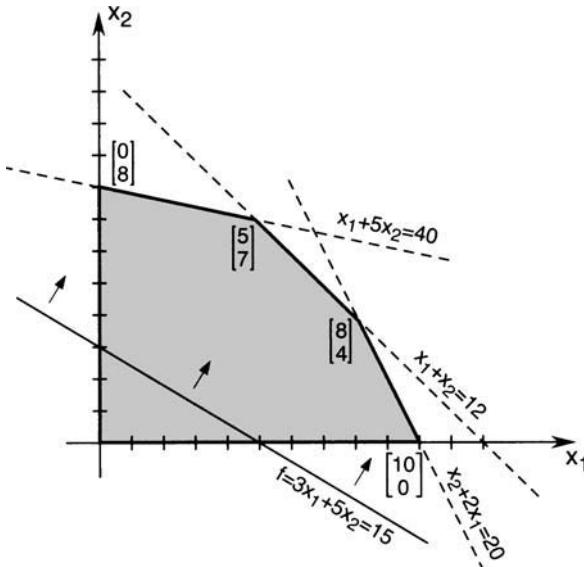
The corresponding basic feasible solution (extreme point) is

$$[10, 0, 30, 0, 2]^\top.$$

For this solution, the objective function value is  $-30$ , which is an improvement relative to the objective function value at the old extreme point.

We now apply the same procedure as above to move to another adjacent extreme point, which hopefully further decreases the value of the objective function. This time, we choose  $\mathbf{a}_2$  to enter the new basis. We have

$$\mathbf{a}_2 = \frac{1}{2}\mathbf{a}_1 + \frac{9}{2}\mathbf{a}_3 + \frac{1}{2}\mathbf{a}_5$$



**Figure 15.11** Graphical solution to the LP problem in Example 15.15.

and

$$\left(10 - \frac{1}{2}\varepsilon_2\right)\mathbf{a}_1 + \varepsilon_2\mathbf{a}_2 + \left(30 - \frac{9}{2}\varepsilon_2\right)\mathbf{a}_3 + \left(2 - \frac{1}{2}\varepsilon_2\right)\mathbf{a}_5 = \mathbf{b}.$$

Substituting  $\varepsilon_2 = 4$ , we obtain

$$8\mathbf{a}_1 + 4\mathbf{a}_2 + 12\mathbf{a}_3 = \mathbf{b}.$$

The solution is  $[8, 4, 12, 0, 0]^\top$  and the corresponding value of the objective function is  $-44$ , which is smaller than the value at the previous extreme point. To complete the example we repeat the procedure once more. This time, we select  $\mathbf{a}_4$  and express it as a combination of the vectors in the previous basis,  $\mathbf{a}_1$ ,  $\mathbf{a}_2$ , and  $\mathbf{a}_3$ :

$$\mathbf{a}_4 = \mathbf{a}_1 - \mathbf{a}_2 + 4\mathbf{a}_3,$$

and hence

$$(8 - \varepsilon_3)\mathbf{a}_1 + (4 + \varepsilon_3)\mathbf{a}_2 + (12 - 4\varepsilon_3)\mathbf{a}_3 + \varepsilon_3\mathbf{a}_4 = \mathbf{b}.$$

The largest permissible value for  $\varepsilon_3$  is  $3$ . The corresponding basic feasible solution is  $[5, 7, 0, 3, 0]^\top$ , with an objective function value of  $-50$ . The solution  $[5, 7, 0, 3, 0]^\top$  turns out to be an optimal solution to our problem in standard form. Hence, the solution to the original problem is  $[5, 7]^\top$ , which we can easily obtain graphically (see Figure 15.11). ■

The technique used in this example for moving from one extreme point to an adjacent extreme point is also used in the simplex method for solving LP problems. The simplex method is essentially a refined method of performing these manipulations.

## EXERCISES

**15.1** Convert the following linear programming problem to *standard form*:

$$\begin{aligned} & \text{maximize} && 2x_1 + x_2 \\ & \text{subject to} && 0 \leq x_1 \leq 2 \\ & && x_1 + x_2 \leq 3 \\ & && x_1 + 2x_2 \leq 5 \\ & && x_2 \geq 0. \end{aligned}$$

**15.2** Consider a discrete-time linear system  $x_{k+1} = ax_k + bu_k$ , where  $u_k$  is the input at time  $k$ ,  $x_k$  is the output at time  $k$ , and  $a, b \in \mathbb{R}$  are system parameters. Given an initial condition  $x_0 = 1$ , consider the problem of minimizing the output  $x_2$  at time 2 subject to the constraint that  $|u_i| \leq 1$ ,  $i = 0, 1$ .

Formulate the problem as a linear programming problem, and convert it into standard form.

**15.3** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && c_1|x_1| + c_2|x_2| + \cdots + c_n|x_n| \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $c_i \neq 0$ ,  $i = 1, \dots, n$ . Convert this problem into an equivalent standard form linear programming problem.

*Hint:* Given any  $x \in \mathbb{R}$ , we can find unique numbers  $x^+, x^- \in \mathbb{R}$ ,  $x^+, x^- \geq 0$ , such that  $|x| = x^+ + x^-$  and  $x = x^+ - x^-$ .

**15.4** Does every linear programming problem in standard form have a nonempty feasible set? If “yes,” provide a proof. If “no,” give a specific example.

Does every linear programming problem in standard form (assuming a nonempty feasible set) have an optimal solution? If “yes,” provide a proof. If “no,” give a specific example.

**15.5** Suppose that a computer supplier has two warehouses, one located in city A and another in city B. The supplier receives orders from two customers, one in city C and another in city D. The customer in city C orders 50 units,

and the customer in city D orders 60 units. The number of units at the warehouse in city A is 70, and the number of units at the warehouse in city B is 80. The cost of shipping each unit from A to C is 1, from A to D is 2, from B to C is 3, and from B to D is 4.

Formulate the problem of deciding how many units from each warehouse should be shipped to each customer to minimize the total shipping cost (assuming that the values of units to be shipped are real numbers). Express the problem as an equivalent standard form linear programming problem.

**15.6** Consider a computer network consisting of six computers,  $A$  through  $F$ . The computers are connected according to the following links, with maximum data rates (in Mbps) shown:  $AC$  (10),  $BC$  (7),  $BF$  (3),  $CD$  (8),  $DE$  (12),  $DF$  (4). For example, “ $AC$  (10)” means that computers  $A$  and  $C$  are connected with a link that supports data rates up to 10 Mbps.

Suppose that  $A$  and  $B$  need to send data to  $E$  and  $F$ , respectively (no other communication is taking place in the network). Any path through the given links above may be used as long as the path has no loop. Also, multiple paths (say from  $A$  to  $E$ ) can be used simultaneously. Link bandwidth can be shared as long as the total data rate through the link does not exceed its maximum (the total data rate through a link is the sum of the data rates of communication in both directions).

For every Mbps of data rate the network can support for transmission from  $A$  to  $E$ , we receive 2 dollars. For every Mbps of data rate the network can support for transmission from  $B$  to  $F$ , we receive 3 dollars. Formulate a linear programming problem to represent the goal of maximizing the total revenue. Then, convert this problem into standard form.

*Hint:* Draw a picture of the network, then label each link with the maximum data rate and the paths that share that link.

**15.7** A cereal manufacturer wishes to produce 1000 pounds of a cereal that contains exactly 10% fiber, 2% fat, and 5% sugar (by weight). The cereal is to be produced by combining four items of raw food material in appropriate proportions. These four items have certain combinations of fiber, fat, and sugar content, and are available at various prices per pound:

Item	1	2	3	4
% fiber	3	8	16	4
% fat	6	46	9	9
% sugar	20	5	4	0
Price/lb	2	4	1	2

The manufacturer wishes to find the amounts of each item to be used to produce the cereal in the least expensive way. Formulate the problem as a linear programming problem. What can you say about the existence of a solution to this problem?

**15.8** Suppose that a wireless broadcast system has  $n$  transmitters. Transmitter  $j$  broadcasts at a power of  $p_j \geq 0$ . There are  $m$  locations where the broadcast is to be received. The path gain from transmitter  $j$  to location  $i$  is  $g_{i,j}$ ; that is, the power of the signal transmitted from transmitter  $j$  received at location  $i$  is  $g_{i,j}p_j$ . The total power received at location  $i$  is the sum of the powers received from all the transmitters. Formulate the problem of finding the minimum sum of the powers transmitted subject to the requirement that the power received at each location is at least  $P$ .

**15.9** Consider the system of equations

$$\begin{bmatrix} 2 & -1 & 2 & -1 & 3 \\ 1 & 2 & 3 & 1 & 0 \\ 1 & 0 & -2 & 0 & -5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 14 \\ 5 \\ -10 \end{bmatrix}.$$

Check if the system has basic solutions. If yes, find all basic solutions.

**15.10** Solve the following linear program graphically:

$$\begin{aligned} &\text{maximize} && 2x_1 + 5x_2 \\ &\text{subject to} && 0 \leq x_1 \leq 4 \\ & && 0 \leq x_2 \leq 6 \\ & && x_1 + x_2 \leq 8. \end{aligned}$$

**15.11** The optimization toolbox in MATLAB provides a function, `linprog`, for solving linear programming problems. Use the function `linprog` to solve the problem in Example 15.5. Use the initial condition **0**.



# CHAPTER 16

---

## SIMPLEX METHOD

---

### 16.1 Solving Linear Equations Using Row Operations

The examples in previous chapters illustrate that solving linear programs involves the solution of systems of linear simultaneous algebraic equations. In this section we describe a method for solving a system of  $n$  linear equations in  $n$  unknowns that we use in subsequent sections. The method uses elementary row operations and corresponding elementary matrices. For a discussion of numerical issues involved in solving a system of simultaneous linear algebraic equations, we refer the reader to [41] and [53].

An elementary row operation on a given matrix is an algebraic manipulation of the matrix that corresponds to one of the following:

1. Interchanging any two rows of the matrix
2. Multiplying one of its rows by a real nonzero number
3. Adding a scalar multiple of one row to another row

An elementary row operation on a matrix is equivalent to premultiplying the matrix by a corresponding *elementary matrix*, which we define next.

**Definition 16.1** We call  $\mathbf{E}$  an *elementary matrix of the first kind* if  $\mathbf{E}$  is obtained from the identity matrix  $\mathbf{I}$  by interchanging any two of its rows. ■

An elementary matrix of the first kind formed from  $\mathbf{I}$  by interchanging the  $i$ th and the  $j$ th rows has the form

$$\mathbf{E} = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & 0 & \cdots & 1 \\ & & & & 1 & \\ & & \vdots & & \ddots & \vdots \\ & & & & & 1 \\ & & & 1 & \cdots & 0 \\ & & & & & 1 \\ & & & & & \ddots \\ & & & & & 1 \end{bmatrix} \quad \begin{array}{l} \leftarrow i\text{th row} \\ \leftarrow j\text{th row} \end{array}$$

Note that  $\mathbf{E}$  is invertible and  $\mathbf{E} = \mathbf{E}^{-1}$ .

**Definition 16.2** We call  $\mathbf{E}$  an *elementary matrix of the second kind* if  $\mathbf{E}$  is obtained from the identity matrix  $\mathbf{I}$  by multiplying one of its rows by a real number  $\alpha \neq 0$ . ■

The elementary matrix of the second kind formed from  $\mathbf{I}$  by multiplying the  $i$ th row by  $\alpha \neq 0$  has the form

$$\mathbf{E} = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & 0 \\ & & 1 & & & \\ & & & \alpha & & \\ & & & & 1 & \\ & 0 & & & & \ddots \\ & & & & & 1 \end{bmatrix} \quad \leftarrow i\text{th row}$$

Note that  $\mathbf{E}$  is invertible and

$$\mathbf{E}^{-1} = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & 0 \\ & & 1 & & & \\ & & & 1/\alpha & & \\ & & & & 1 & \\ 0 & & & & & \ddots \\ & & & & & & 1 \end{bmatrix} \quad \leftarrow \text{ith row}$$

**Definition 16.3** We call  $\mathbf{E}$  an *elementary matrix of the third kind* if  $\mathbf{E}$  is obtained from the identity matrix  $\mathbf{I}$  by adding  $\beta$  times one row to another row of  $\mathbf{I}$ . ■

An elementary matrix of the third kind obtained from  $\mathbf{I}$  by adding  $\beta$  times the  $j$ th row to the  $i$ th row has the form

$$\mathbf{E} = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & 0 \\ & & 1 & \cdots & \beta & \\ & & & \ddots & \vdots & \\ & & & & 1 & \\ 0 & & & & & \ddots \\ & & & & & & 1 \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{ith row} \\ \leftarrow \text{jth row} \end{array}$$

Observe that  $\mathbf{E}$  is the identity matrix with an extra  $\beta$  in the  $(i,j)$ th location. Note that  $\mathbf{E}$  is invertible and

$$\mathbf{E}^{-1} = \begin{bmatrix} 1 & & & & & & \\ & \ddots & & & & & 0 \\ & & 1 & \cdots & -\beta & & \\ & & & \ddots & \vdots & & \\ & & & & 1 & & \\ 0 & & & & & \ddots & \\ & & & & & & 1 \end{bmatrix} \quad \begin{array}{l} \leftarrow \text{ith row} \\ \leftarrow \text{jth row} \end{array}$$

**Definition 16.4** An *elementary row operation* (of first, second, or third kind) on a given matrix is a premultiplication of the given matrix by a corresponding elementary matrix of the respective kind. ■

Because elementary matrices are invertible, we can define the corresponding inverse elementary row operations.

Consider a system of  $n$  linear equations in  $n$  unknowns  $x_1, x_2, \dots, x_n$  with right-hand sides  $b_1, b_2, \dots, b_n$ . In matrix form this system may be written as

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}.$$

If  $\mathbf{A}$  is invertible, then

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}.$$

Thus, the problem of solving the system of equations  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , with  $\mathbf{A} \in \mathbb{R}^{n \times n}$  invertible, is related to the problem of computing  $\mathbf{A}^{-1}$ . We now show that  $\mathbf{A}^{-1}$  can be computed effectively using elementary row operations. In particular, we prove the following theorem.

**Theorem 16.1** *Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a given matrix. Then,  $\mathbf{A}$  is nonsingular (invertible) if and only if there exist elementary matrices  $\mathbf{E}_i$ ,  $i = 1, \dots, t$ , such that*

$$\mathbf{E}_t \cdots \mathbf{E}_2 \mathbf{E}_1 \mathbf{A} = \mathbf{I}.$$

□

*Proof.*  $\Rightarrow$ : If  $\mathbf{A}$  is nonsingular, then its first column must have at least one nonzero element, say  $a_{j1} \neq 0$ . Premultiplying  $\mathbf{A}$  by an elementary matrix of the first kind of the form

$$\mathbf{E}_1 = \left[ \begin{array}{cccccc} 0 & & & 1 & & & \\ & 1 & & & & & \\ & & \ddots & & \vdots & & \\ & & & 1 & & & \\ 1 & \dots & 0 & & 1 & & \\ & & & & & \ddots & \\ & & & & & & 1 \end{array} \right] \quad \text{← } j\text{th row}$$

brings the nonzero element  $a_{j1}$  to the location  $(1, 1)$ . Hence, in the matrix  $\mathbf{E}_1 \mathbf{A}$ , the element  $a_{11} \neq 0$ . Note that since  $\mathbf{E}_1$  is nonsingular,  $\mathbf{E}_1 \mathbf{A}$  is also nonsingular.

Next, we premultiply  $\mathbf{E}_1\mathbf{A}$  by an elementary matrix of the second kind of the form

$$\mathbf{E}_2 = \begin{bmatrix} 1/a_{11} & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}.$$

The result of this operation is the matrix  $\mathbf{E}_2\mathbf{E}_1\mathbf{A}$  with unity in the location  $(1, 1)$ . We next apply a sequence of elementary row operations of the third kind on the matrix  $\mathbf{E}_2\mathbf{E}_1\mathbf{A}$ . Specifically, we premultiply  $\mathbf{E}_2\mathbf{E}_1\mathbf{A}$  by  $n - 1$  elementary matrices of the form

$$\mathbf{E}_3 = \begin{bmatrix} 1 & & & \\ -a_{21} & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix}, \dots, \mathbf{E}_r = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ \vdots & & 1 & \\ -a_{n1} & & & \ddots \\ & & & & 1 \end{bmatrix},$$

where  $r = 2 + n - 1 = n + 1$ . The result of these operations is the nonsingular matrix

$$\mathbf{E}_r\mathbf{E}_{r-1}\cdots\mathbf{E}_2\mathbf{E}_1\mathbf{A} = \begin{bmatrix} 1 & \bar{a}_{12} & \cdots & \bar{a}_{1n} \\ 0 & \bar{a}_{22} & \cdots & \bar{a}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \bar{a}_{n2} & \cdots & \bar{a}_{nn} \end{bmatrix}.$$

Because the matrix  $\mathbf{E}_r\cdots\mathbf{E}_1\mathbf{A}$  is nonsingular, its submatrix

$$\begin{bmatrix} \bar{a}_{22} & \cdots & \bar{a}_{2n} \\ \vdots & & \vdots \\ \bar{a}_{n2} & \cdots & \bar{a}_{nn} \end{bmatrix}$$

must also be nonsingular. This implies that there is a nonzero element  $\bar{a}_{j2}$ , where  $2 \leq j \leq n$ . Using an elementary operation of the first kind, we bring this element to the location  $(2, 2)$ . Thus, in the matrix

$$\mathbf{E}_{r+1}\mathbf{E}_r\cdots\mathbf{E}_1\mathbf{A}$$

the  $(2, 2)$ th element is nonzero. Premultiplying the matrix by an elementary matrix of the second kind yields the matrix

$$\mathbf{E}_{r+2}\mathbf{E}_{r+1}\cdots\mathbf{E}_1\mathbf{A},$$

in which the element in the location  $(2, 2)$  is unity. As before, we premultiply this matrix by  $n - 1$  elementary row operations of the third kind, to get a

matrix of the form

$$\mathbf{E}_s \cdots \mathbf{E}_r \cdots \mathbf{E}_1 \mathbf{A} = \begin{bmatrix} 1 & 0 & \tilde{a}_{13} & \cdots & \tilde{a}_{1n} \\ 0 & 1 & \tilde{a}_{23} & \cdots & \tilde{a}_{2n} \\ 0 & 0 & \tilde{a}_{33} & \cdots & \tilde{a}_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \tilde{a}_{n3} & \cdots & \tilde{a}_{nn} \end{bmatrix},$$

where  $s = r + 2 + n - 1 = 2(n + 1)$ . This matrix is nonsingular. Hence, there is a nonzero element  $\tilde{a}_{j3}$ ,  $3 \leq j \leq n$ . Proceeding in a similar fashion as before, we obtain

$$\mathbf{E}_t \cdots \mathbf{E}_s \cdots \mathbf{E}_r \cdots \mathbf{E}_1 \mathbf{A} = \mathbf{I},$$

where  $t = n(n + 1)$ .

$\Leftarrow$ : If there exist elementary matrices  $\mathbf{E}_1, \dots, \mathbf{E}_t$  such that

$$\mathbf{E}_t \cdots \mathbf{E}_1 \mathbf{A} = \mathbf{I},$$

then clearly  $\mathbf{A}$  is invertible, with

$$\mathbf{A}^{-1} = \mathbf{E}_t \cdots \mathbf{E}_1.$$

■

Theorem 16.1 suggests the following procedure for finding  $\mathbf{A}^{-1}$ , if it exists. We first form the matrix

$$[\mathbf{A}, \mathbf{I}].$$

We then apply elementary row operations to  $[\mathbf{A}, \mathbf{I}]$  so that  $\mathbf{A}$  is transformed into  $\mathbf{I}$ ; that is, we obtain

$$\mathbf{E}_t \cdots \mathbf{E}_1 [\mathbf{A}, \mathbf{I}] = [\mathbf{I}, \mathbf{B}].$$

It then follows that

$$\mathbf{B} = \mathbf{E}_t \cdots \mathbf{E}_1 = \mathbf{A}^{-1}.$$

**Example 16.1** Let

$$\mathbf{A} = \begin{bmatrix} 2 & 5 & 10 & 0 \\ 1 & 1 & 1 & 0 \\ -2 & -10 & -30 & 1 \\ -1 & -2 & -3 & 0 \end{bmatrix}.$$

Find  $\mathbf{A}^{-1}$ .

We form the matrix

$$[\mathbf{A}, \mathbf{I}] = \begin{bmatrix} 2 & 5 & 10 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ -2 & -10 & -30 & 1 & 0 & 0 & 1 & 0 \\ -1 & -2 & -3 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

and perform row operations on this matrix. Applying row operations of the first and third kinds yields

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 3 & 8 & 0 & 1 & -2 & 0 & 0 \\ 0 & -8 & -28 & 1 & 0 & 2 & 1 & 0 \\ 0 & -1 & -2 & 0 & 0 & 1 & 0 & 1 \end{bmatrix}.$$

We then interchange the second and fourth rows and apply elementary row operations of the second and third kinds to get

$$\begin{bmatrix} 1 & 0 & -1 & 0 & 0 & 2 & 0 & 1 \\ 0 & 1 & 2 & 0 & 0 & -1 & 0 & -1 \\ 0 & 0 & 2 & 0 & 1 & 1 & 0 & 3 \\ 0 & 0 & -12 & 1 & 0 & -6 & 1 & -8 \end{bmatrix}.$$

Now multiply the third row by  $1/2$  and then perform a sequence of elementary operations of the third kind to obtain

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \frac{1}{2} & \frac{5}{2} & 0 & \frac{5}{2} \\ 0 & 1 & 0 & 0 & -1 & -2 & 0 & -4 \\ 0 & 0 & 1 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & \frac{3}{2} \\ 0 & 0 & 0 & 1 & 6 & 0 & 1 & 10 \end{bmatrix}.$$

Hence,

$$\mathbf{A}^{-1} = \begin{bmatrix} \frac{1}{2} & \frac{5}{2} & 0 & \frac{5}{2} \\ -1 & -2 & 0 & -4 \\ \frac{1}{2} & \frac{1}{2} & 0 & \frac{3}{2} \\ 6 & 0 & 1 & 10 \end{bmatrix}.$$

■

We now return to the general problem of solving the system of equations  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . If  $\mathbf{A}^{-1}$  exists, then the solution is  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ . However, we do not need an explicit expression for  $\mathbf{A}^{-1}$  to find the solution. Indeed, let  $\mathbf{A}^{-1}$  be expressed as a product of elementary matrices

$$\mathbf{A}^{-1} = \mathbf{E}_t \mathbf{E}_{t-1} \cdots \mathbf{E}_1.$$

Thus,

$$\mathbf{E}_t \cdots \mathbf{E}_1 \mathbf{A} \mathbf{x} = \mathbf{E}_t \cdots \mathbf{E}_1 \mathbf{b}$$

and hence

$$\mathbf{x} = \mathbf{E}_t \cdots \mathbf{E}_1 \mathbf{b}.$$

The discussion above leads to the following procedure for solving the system  $\mathbf{Ax} = \mathbf{b}$ . Form an augmented matrix

$$[\mathbf{A}, \mathbf{b}].$$

Then, perform a sequence of row elementary operations on this augmented matrix until we obtain

$$[\mathbf{I}, \tilde{\mathbf{b}}].$$

From the above we have that if  $\mathbf{x}$  is a solution to  $\mathbf{Ax} = \mathbf{b}$ , then it is also a solution to  $\mathbf{EAx} = \mathbf{Eb}$ , where  $\mathbf{E} = \mathbf{E}_t \cdots \mathbf{E}_1$  represents a sequence of elementary row operations. Because  $\mathbf{EA} = \mathbf{I}$ , and  $\mathbf{Eb} = \tilde{\mathbf{b}}$ , it follows that  $\mathbf{x} = \tilde{\mathbf{b}}$  is the solution to  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$  invertible.

Suppose now that  $\mathbf{A} \in \mathbb{R}^{m \times n}$  where  $m < n$ , and  $\text{rank } \mathbf{A} = m$ . Then,  $\mathbf{A}$  is not a square matrix. Clearly, in this case the system of equations  $\mathbf{Ax} = \mathbf{b}$  has infinitely many solutions. Without loss of generality, we can assume that the first  $m$  columns of  $\mathbf{A}$  are linearly independent. Then, if we perform a sequence of elementary row operations on the augmented matrix  $[\mathbf{A}, \mathbf{b}]$  as before, we obtain

$$[\mathbf{I}, \mathbf{D}, \tilde{\mathbf{b}}],$$

where  $\mathbf{D}$  is an  $m \times (n - m)$  matrix. Let  $\mathbf{x} \in \mathbb{R}^n$  be a solution to  $\mathbf{Ax} = \mathbf{b}$  and write  $\mathbf{x} = [\mathbf{x}_B^\top, \mathbf{x}_D^\top]^\top$ , where  $\mathbf{x}_B \in \mathbb{R}^m$ ,  $\mathbf{x}_D \in \mathbb{R}^{(n-m)}$ . Then,  $[\mathbf{I}, \mathbf{D}]\mathbf{x} = \tilde{\mathbf{b}}$ , which we can rewrite as  $\mathbf{x}_B + \mathbf{D}\mathbf{x}_D = \tilde{\mathbf{b}}$ , or  $\mathbf{x}_B = \tilde{\mathbf{b}} - \mathbf{D}\mathbf{x}_D$ . Note that for an arbitrary  $\mathbf{x}_D \in \mathbb{R}^{(n-m)}$ , if  $\mathbf{x}_B = \tilde{\mathbf{b}} - \mathbf{D}\mathbf{x}_D$ , then the resulting vector  $\mathbf{x} = [\mathbf{x}_B^\top, \mathbf{x}_D^\top]^\top$  is a solution to  $\mathbf{Ax} = \mathbf{b}$ . In particular,  $[\tilde{\mathbf{b}}^\top, \mathbf{0}^\top]^\top$  is a solution to  $\mathbf{Ax} = \mathbf{b}$ . We often refer to the basic solution  $[\tilde{\mathbf{b}}^\top, \mathbf{0}^\top]^\top$  as a *particular solution* to  $\mathbf{Ax} = \mathbf{b}$ . Note that  $[-(\mathbf{D}\mathbf{x}_D)^\top, \mathbf{x}_D^\top]^\top$  is a solution to  $\mathbf{Ax} = \mathbf{0}$ . Any solution to  $\mathbf{Ax} = \mathbf{b}$  has the form

$$\mathbf{x} = \begin{bmatrix} \tilde{\mathbf{b}} \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} -\mathbf{D}\mathbf{x}_D \\ \mathbf{x}_D \end{bmatrix}$$

for some  $\mathbf{x}_D \in \mathbb{R}^{(n-m)}$ .

## 16.2 The Canonical Augmented Matrix

Consider the system of simultaneous linear equations  $\mathbf{Ax} = \mathbf{b}$ ,  $\text{rank } \mathbf{A} = m$ . Using a sequence of elementary row operations and reordering the variables

if necessary, we transform the system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  into the following *canonical form*:

$$\begin{aligned} x_1 &+ y_{1m+1}x_{m+1} + \cdots + y_{1n}x_n = y_{10} \\ x_2 &+ y_{2m+1}x_{m+1} + \cdots + y_{2n}x_n = y_{20} \\ &\vdots \\ x_m &+ y_{mm+1}x_{m+1} + \cdots + y_{mn}x_n = y_{m0}. \end{aligned}$$

This can be represented in matrix notation as

$$[\mathbf{I}_m, \mathbf{Y}_{m,n-m}] \mathbf{x} = \mathbf{y}_0.$$

Formally, we define the canonical form as follows.

**Definition 16.5** A system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is said to be in *canonical form* if among the  $n$  variables there are  $m$  variables with the property that each appears in only one equation, and its coefficient in that equation is unity. ■

A system is in canonical form if by some reordering of the equations and the variables it takes the form  $[\mathbf{I}_m, \mathbf{Y}_{m,n-m}] \mathbf{x} = \mathbf{y}_0$ . If a system of equations  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is not in canonical form, we can transform the system into canonical form by a sequence of elementary row operations. The system in canonical form has the same solution as the original system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and is called the *canonical representation* of the system with respect to the basis  $\mathbf{a}_1, \dots, \mathbf{a}_m$ . There are, in general, many canonical representations of a given system, depending on which columns of  $\mathbf{A}$  we transform into the columns of  $\mathbf{I}_m$  (i.e., basic columns). We call the augmented matrix  $[\mathbf{I}_m, \mathbf{Y}_{m,n-m}, \mathbf{y}_0]$  of the canonical representation of a given system the *canonical augmented matrix* of the system with respect to the basis  $\mathbf{a}_1, \dots, \mathbf{a}_m$ . Of course, there may be many canonical augmented matrices of a given system, depending on which columns of  $\mathbf{A}$  are chosen as basic columns.

The variables corresponding to basic columns in a canonical representation of a given system are the basic variables, whereas the other variables are the nonbasic variables. In particular, in the canonical representation  $[\mathbf{I}_m, \mathbf{Y}_{m,n-m}] \mathbf{x} = \mathbf{y}_0$  of a given system, the variables  $x_1, \dots, x_m$  are the basic variables and the other variables are the nonbasic variables. Note that in general the basic variables need not be the first  $m$  variables. However, in the following discussion we assume, for convenience and without loss of generality, that the basic variables are indeed the first  $m$  variables in the system.

Having done so, the corresponding basic solution is

$$x_1 = y_{10},$$

$$\vdots$$

$$x_m = y_{m0},$$

$$x_{m+1} = 0,$$

$$\vdots$$

$$x_n = 0;$$

that is,

$$\mathbf{x} = \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{0} \end{bmatrix}.$$

Given a system of equations  $\mathbf{Ax} = \mathbf{b}$ , consider the associated canonical augmented matrix

$$[\mathbf{I}_m, \mathbf{Y}_{m,n-m}, \mathbf{y}_0] = \begin{bmatrix} 1 & 0 & \cdots & 0 & y_{1m+1} & \cdots & y_{1n} & y_{10} \\ 0 & 1 & \cdots & 0 & y_{2m+1} & \cdots & y_{2n} & y_{20} \\ \vdots & \vdots & \ddots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & y_{mm+1} & \cdots & y_{mn} & y_{m0} \end{bmatrix}.$$

From the arguments above we conclude that

$$\mathbf{b} = y_{10}\mathbf{a}_1 + y_{20}\mathbf{a}_2 + \cdots + y_{m0}\mathbf{a}_m.$$

In other words, the entries in the last column of the canonical augmented matrix are the coordinates of the vector  $\mathbf{b}$  with respect to the basis  $\{\mathbf{a}_1, \dots, \mathbf{a}_m\}$ . The entries of all the other columns of the canonical augmented matrix have a similar interpretation. Specifically, the entries of the  $j$ th column of the canonical augmented matrix,  $j = 1, \dots, n$ , are the coordinates of  $\mathbf{a}_j$  with respect to the basis  $\{\mathbf{a}_1, \dots, \mathbf{a}_m\}$ . To see this, note that the first  $m$  columns of the augmented matrix form a basis (the standard basis). Every other vector in the augmented matrix can be expressed as a linear combination of these basis vectors by reading the coefficients down the corresponding column. Specifically, let  $\mathbf{a}'_i$ ,  $i = 1, \dots, n+1$ , be the  $i$ th column in the augmented matrix above. Clearly, since  $\mathbf{a}'_1, \dots, \mathbf{a}'_m$  form the standard basis, then for  $m < j \leq n$ ,

$$\mathbf{a}'_j = y_{1j}\mathbf{a}'_1 + y_{2j}\mathbf{a}'_2 + \cdots + y_{mj}\mathbf{a}'_m.$$

Let  $\mathbf{a}_i$ ,  $i = 1, \dots, n$ , be the  $i$ th column of  $\mathbf{A}$ , and  $\mathbf{a}_{n+1} = \mathbf{b}$ . Now,  $\mathbf{a}'_i = \mathbf{E}\mathbf{a}_i$ ,  $i = 1, \dots, n+1$ , where  $\mathbf{E}$  is a nonsingular matrix that represents the elementary row operations needed to transform  $[\mathbf{A}, \mathbf{b}]$  into  $[\mathbf{I}_m, \mathbf{Y}_{m,n-m}, \mathbf{y}_0]$ . Therefore, for  $m < j \leq n$ , we also have

$$\mathbf{a}_j = y_{1j}\mathbf{a}_1 + y_{2j}\mathbf{a}_2 + \cdots + y_{mj}\mathbf{a}_m.$$

### 16.3 Updating the Augmented Matrix

To summarize Section 16.2, the canonical augmented matrix of a given system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  specifies the representations of the columns  $\mathbf{a}_j$ ,  $m < j \leq n$ , in terms of the basic columns  $\mathbf{a}_1, \dots, \mathbf{a}_m$ . Thus, the elements of the  $j$ th column of the canonical augmented matrix are the coordinates of the vector  $\mathbf{a}_j$  with respect to the basis  $\mathbf{a}_1, \dots, \mathbf{a}_m$ . The coordinates of  $\mathbf{b}$  are given in the last column.

Suppose that we are given the canonical representation of a system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . We now consider the following question: If we replace a basic variable by a nonbasic variable, what is the new canonical representation corresponding to the new set of basic variables? Specifically, suppose that we wish to replace the basis vector  $\mathbf{a}_p$ ,  $1 \leq p \leq m$ , by the vector  $\mathbf{a}_q$ ,  $m < q \leq n$ . Provided that the first  $m$  vectors with  $\mathbf{a}_p$  replaced by  $\mathbf{a}_q$  are linearly independent, these vectors constitute a basis and every vector can be expressed as a linear combination of the new basic columns.

Let us now find the coordinates of the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_n$  with respect to the new basis. These coordinates form the entries of the canonical augmented matrix of the system with respect to the new basis. In terms of the old basis, we can express  $\mathbf{a}_q$  as

$$\mathbf{a}_q = \sum_{i=1}^m y_{iq} \mathbf{a}_i = \sum_{\substack{i=1 \\ i \neq p}}^m y_{iq} \mathbf{a}_i + y_{pq} \mathbf{a}_p.$$

Note that the set of vectors  $\{\mathbf{a}_1, \dots, \mathbf{a}_{p-1}, \mathbf{a}_q, \mathbf{a}_{p+1}, \dots, \mathbf{a}_m\}$  is linearly independent if and only if  $y_{pq} \neq 0$ . Solving the equation above for  $\mathbf{a}_p$ , we get

$$\mathbf{a}_p = \frac{1}{y_{pq}} \mathbf{a}_q - \sum_{\substack{i=1 \\ i \neq p}}^m \frac{y_{iq}}{y_{pq}} \mathbf{a}_i.$$

Recall that in terms of the old augmented matrix, any vector  $\mathbf{a}_j$ ,  $m < j \leq n$ , can be expressed as

$$\mathbf{a}_j = y_{1j} \mathbf{a}_1 + y_{2j} \mathbf{a}_2 + \cdots + y_{mj} \mathbf{a}_m.$$

Combining the last two equations yields

$$\mathbf{a}_j = \sum_{\substack{i=1 \\ i \neq p}}^m \left( y_{ij} - \frac{y_{pj}}{y_{pq}} y_{iq} \right) \mathbf{a}_i + \frac{y_{pj}}{y_{pq}} \mathbf{a}_q.$$

Denoting the entries of the new augmented matrix by  $y'_{ij}$ , we obtain

$$y'_{ij} = y_{ij} - \frac{y_{pj}}{y_{pq}} y_{iq}, \quad i \neq p,$$

$$y'_{pj} = \frac{y_{pj}}{y_{pq}}.$$

Therefore, the entries of the new canonical augmented matrix can be obtained from the entries of the old canonical augmented matrix via the formulas above. These equations are often called the *pivot equations*, and  $y_{pq}$ , the *pivot element*.

We refer to the operation on a given matrix by the formulas above as *pivoting about the  $(p, q)$ th element*. Note that pivoting about the  $(p, q)$ th element results in a matrix whose  $q$ th column has all zero entries, except the  $(p, q)$ th entry, which is unity. The pivoting operation can be accomplished via a sequence of elementary row operations, as was done in the proof of Theorem 16.1.

## 16.4 The Simplex Algorithm

The essence of the simplex algorithm is to move from one basic feasible solution to another until an optimal basic feasible solution is found. The canonical augmented matrix discussed in Section 16.3 plays a central role in the simplex algorithm.

Suppose that we are given the basic feasible solution

$$\mathbf{x} = [x_1, \dots, x_m, 0, \dots, 0]^\top, \quad x_i \geq 0, \quad i = 1, \dots, m$$

or equivalently

$$x_1 \mathbf{a}_1 + \cdots + x_m \mathbf{a}_m = \mathbf{b}.$$

In Section 16.3 we saw how to update the canonical augmented matrix if we wish to replace a basic column by a nonbasic column, that is, if we wish to change from one basis to another by replacing a single basic column. The values of the basic variables in a basic solution corresponding to a given basis are given in the last column of the canonical augmented matrix with respect to that basis; that is,  $x_i = y_{i0}$ ,  $i = 1, \dots, m$ . Basic solutions are not necessarily feasible—the values of the basic variables may be negative. In the simplex method we want to move from one basic feasible solution to another. This means that we want to change basic columns in such a way that the last column of the canonical augmented matrix remains nonnegative. In this section we discuss a systematic method for doing this.

In the remainder of this chapter we assume that every basic feasible solution of

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0} \end{aligned}$$

is a nondegenerate basic feasible solution. We make this assumption primarily for convenience—all arguments can be extended to include degeneracy.

Let us start with the basic columns  $\mathbf{a}_1, \dots, \mathbf{a}_m$ , and assume that the corresponding basic solution  $\mathbf{x} = [y_{10}, \dots, y_{m0}, 0, \dots, 0]^\top$  is feasible; that is, the

entries  $y_{i0}$ ,  $i = 1, \dots, m$ , in the last column of the canonical augmented matrix are positive. Suppose that we now decide to make the vector  $\mathbf{a}_q$ ,  $q > m$ , a basic column. We first represent  $\mathbf{a}_q$  in terms of the current basis as

$$\mathbf{a}_q = y_{1q}\mathbf{a}_1 + y_{2q}\mathbf{a}_2 + \cdots + y_{mq}\mathbf{a}_m.$$

Multiplying the above by  $\varepsilon > 0$  yields

$$\varepsilon\mathbf{a}_q = \varepsilon y_{1q}\mathbf{a}_1 + \varepsilon y_{2q}\mathbf{a}_2 + \cdots + \varepsilon y_{mq}\mathbf{a}_m.$$

We combine this equation with

$$y_{10}\mathbf{a}_1 + \cdots + y_{m0}\mathbf{a}_m = \mathbf{b}$$

to get

$$(y_{10} - \varepsilon y_{1q})\mathbf{a}_1 + (y_{20} - \varepsilon y_{2q})\mathbf{a}_2 + \cdots + (y_{m0} - \varepsilon y_{mq})\mathbf{a}_m + \varepsilon\mathbf{a}_q = \mathbf{b}.$$

Note that the vector

$$\begin{bmatrix} y_{10} - \varepsilon y_{1q} \\ \vdots \\ y_{m0} - \varepsilon y_{mq} \\ 0 \\ \vdots \\ \varepsilon \\ \vdots \\ 0 \end{bmatrix},$$

where  $\varepsilon$  appears in the  $q$ th position, is a solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . If  $\varepsilon = 0$ , then we obtain the old basic feasible solution. As  $\varepsilon$  is increased from zero, the  $q$ th component of the vector above increases. All other entries of this vector will increase or decrease linearly as  $\varepsilon$  is increased, depending on whether the corresponding  $y_{iq}$  is negative or positive. For small enough  $\varepsilon$ , we have a feasible but nonbasic solution. If any of the components decreases as  $\varepsilon$  increases, we choose  $\varepsilon$  to be the smallest value where one (or more) of the components vanishes. That is,

$$\varepsilon = \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}.$$

With this choice of  $\varepsilon$  we have a new basic feasible solution, with the vector  $\mathbf{a}_q$  replacing  $\mathbf{a}_p$ , where  $p$  corresponds to the minimizing index  $p = \arg \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}$ . So, we now have a new basis  $\mathbf{a}_1, \dots, \mathbf{a}_{p-1}, \mathbf{a}_{p+1}, \dots, \mathbf{a}_m, \mathbf{a}_q$ . As we can see,  $\mathbf{a}_p$  was replaced by  $\mathbf{a}_q$  in the new basis. We say that  $\mathbf{a}_q$  enters the basis and  $\mathbf{a}_p$  leaves the basis. If the minimum in  $\min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}$  is achieved by more than a single index, then the new solution is degenerate and any of the zero components can

be regarded as the component corresponding to the basic column that leaves the basis. If none of the  $y_{iq}$  are positive, then all components in the vector  $[y_{10} - \varepsilon y_{1q}, \dots, y_{m0} - \varepsilon y_{mq}, 0, \dots, \varepsilon, \dots, 0]^\top$  increase (or remain constant) as  $\varepsilon$  is increased, and no new basic feasible solution is obtained, no matter how large we make  $\varepsilon$ . In this case there are feasible solutions having arbitrarily large components, which means that the set  $\Omega$  of feasible solutions is unbounded.

So far, we have discussed how to change from one basis to another, while preserving feasibility of the corresponding basic solution, assuming that we have already chosen a nonbasic column to enter the basis. To complete our development of the simplex method, we need to consider two more issues. The first issue concerns the choice of which nonbasic column should enter the basis. The second issue is to find a stopping criterion, that is, a way to determine if a basic feasible solution is optimal or is not. To this end, suppose that we have found a basic feasible solution. The main idea of the simplex method is to move from one basic feasible solution (extreme point of the set  $\Omega$ ) to another basic feasible solution at which the value of the objective function is smaller. Because there is only a finite number of extreme points of the feasible set, the optimal point will be reached after a finite number of steps.

We already know how to move from one extreme point of the set  $\Omega$  to a neighboring one by updating the canonical augmented matrix. To see which neighboring solution we should move to and when to stop moving, consider the following basic feasible solution:

$$[\mathbf{x}_B^\top, \mathbf{0}^\top]^\top = [y_{10}, \dots, y_{m0}, 0, \dots, 0]^\top$$

together with the corresponding canonical augmented matrix, having an identity matrix appearing in the first  $m$  columns. The value of the objective function for any solution  $\mathbf{x}$  is

$$z = c_1 x_1 + c_2 x_2 + \cdots + c_n x_n.$$

For our basic solution, the value of the objective function is

$$z = z_0 = \mathbf{c}_B^\top \mathbf{x}_B = c_1 y_{10} + \cdots + c_m y_{m0},$$

where

$$\mathbf{c}_B^\top = [c_1, c_2, \dots, c_m].$$

To see how the value of the objective function changes when we move from one basic feasible solution to another, suppose that we choose the  $q$ th column,  $m < q \leq n$ , to enter the basis. To update the canonical augmented matrix, let  $p = \arg \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}$  and  $\varepsilon = y_{p0}/y_{pq}$ . The new basic feasible

solution is

$$\begin{bmatrix} y_{10} - \varepsilon y_{1q} \\ \vdots \\ y_{m0} - \varepsilon y_{mq} \\ 0 \\ \vdots \\ \varepsilon \\ \vdots \\ 0 \end{bmatrix}.$$

Note that the single  $\varepsilon$  appears in the  $q$ th component, whereas the  $p$ th component is zero. Observe that we could have arrived at the basic feasible solution above simply by updating the canonical augmented matrix using the pivot equations from the previous Section 16.3:

$$\begin{aligned} y'_{ij} &= y_{ij} - \frac{y_{pj}}{y_{pq}} y_{iq}, \quad i \neq p, \\ y'_{pj} &= \frac{y_{pj}}{y_{pq}}, \end{aligned}$$

where the  $q$ th column enters the basis and the  $p$ th column leaves [i.e., we pivot about the  $(p, q)$ th element]. The values of the basic variables are entries in the last column of the updated canonical augmented matrix.

The cost for this new basic feasible solution is

$$\begin{aligned} z &= c_1(y_{10} - y_{1q}\varepsilon) + \cdots + c_m(y_{m0} - y_{mq}\varepsilon) + c_q\varepsilon \\ &= z_0 + [c_q - (c_1y_{1q} + \cdots + c_my_{mq})]\varepsilon, \end{aligned}$$

where  $z_0 = c_1y_{10} + \cdots + c_my_{m0}$ . Let

$$z_q = c_1y_{1q} + \cdots + c_my_{mq}.$$

Then,

$$z = z_0 + (c_q - z_q)\varepsilon.$$

Thus, if

$$z - z_0 = (c_q - z_q)\varepsilon < 0,$$

then the objective function value at the new basic feasible solution above is smaller than the objective function value at the original solution (i.e.,  $z < z_0$ ). Therefore, if  $c_q - z_q < 0$ , then the new basic feasible solution with  $a_q$  entering the basis has a lower objective function value.

On the other hand, if the given basic feasible solution is such that for all  $q = m+1, \dots, n$ ,

$$c_q - z_q \geq 0,$$

then we can show that this solution is in fact an optimal solution. To show this, recall from Section 16.1 that any solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$  can be represented as

$$\mathbf{x} = \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} -\mathbf{Y}_{m,n-m} \mathbf{x}_D \\ \mathbf{x}_D \end{bmatrix}$$

for some  $\mathbf{x}_D = [x_{m+1}, \dots, x_n]^\top \in \mathbb{R}^{(n-m)}$ . Using manipulations similar to the above, we obtain

$$\mathbf{c}^\top \mathbf{x} = z_0 + \sum_{i=m+1}^n (c_i - z_i)x_i,$$

where  $z_i = c_1 y_{1i} + \dots + c_m y_{mi}$ ,  $i = m+1, \dots, n$ . For a feasible solution we have  $x_i \geq 0$ ,  $i = 1, \dots, n$ . Therefore, if  $c_i - z_i \geq 0$  for all  $i = m+1, \dots, n$ , then any feasible solution  $\mathbf{x}$  will have objective function value  $\mathbf{c}^\top \mathbf{x}$  no smaller than  $z_0$ .

Let  $r_i = 0$  for  $i = 1, \dots, m$  and  $r_i = c_i - z_i$  for  $i = m+1, \dots, n$ . We call  $r_i$  the *i*th *reduced cost coefficient* or *relative cost coefficient*. Note that the reduced cost coefficients corresponding to basic variables are zero.

We summarize the discussion above with the following result.

**Theorem 16.2** *A basic feasible solution is optimal if and only if the corresponding reduced cost coefficients are all nonnegative.*  $\square$

At this point we have all the necessary steps for the simplex algorithm.

### Simplex Algorithm

1. Form a canonical augmented matrix corresponding to an initial basic feasible solution.
2. Calculate the reduced cost coefficients corresponding to the nonbasic variables.
3. If  $r_j \geq 0$  for all  $j$ , stop—the current basic feasible solution is optimal.
4. Select a  $q$  such that  $r_q < 0$ .
5. If no  $y_{iq} > 0$ , stop—the problem is unbounded; else, calculate  $p = \arg \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}$ . (If more than one index  $i$  minimizes  $y_{i0}/y_{iq}$ , we let  $p$  be the smallest such index.)
6. Update the canonical augmented matrix by pivoting about the  $(p, q)$ th element.
7. Go to step 2.

We state the following result for the simplex algorithm, which we have already proved in the foregoing discussion.

**Theorem 16.3** Suppose that we have an LP problem in standard form that has an optimal feasible solution. If the simplex method applied to this problem terminates and the reduced cost coefficients in the last step are all nonnegative, then the resulting basic feasible solution is optimal.  $\square$

**Example 16.2** Consider the following linear program (see also Exercise 15.10):

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 \\ & \text{subject to} && x_1 \leq 4 \\ & && x_2 \leq 6 \\ & && x_1 + x_2 \leq 8 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

We solve this problem using the simplex method.

Introducing slack variables, we transform the problem into standard form:

$$\begin{aligned} & \text{minimize} && -2x_1 - 5x_2 - 0x_3 - 0x_4 - 0x_5 \\ & \text{subject to} && x_1 + x_3 = 4 \\ & && x_2 + x_4 = 6 \\ & && x_1 + x_2 + x_5 = 8 \\ & && x_1, x_2, x_3, x_4, x_5 \geq 0. \end{aligned}$$

The starting canonical augmented matrix for this problem is

$$\begin{array}{ccccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 1 & 0 & 1 & 0 & 0 & 4 \\ 0 & 1 & 0 & 1 & 0 & 6 \\ 1 & 1 & 0 & 0 & 1 & 8 \end{array}$$

Observe that the columns forming the identity matrix in the canonical augmented matrix above do not appear at the beginning. We could rearrange the augmented matrix so that the identity matrix would appear first. However, this is not essential from the computational point of view.

The starting basic feasible solution to the problem in standard form is

$$\mathbf{x} = [0, 0, 4, 6, 8]^T.$$

The columns  $\mathbf{a}_3$ ,  $\mathbf{a}_4$ , and  $\mathbf{a}_5$  corresponding to  $x_3$ ,  $x_4$ , and  $x_5$  are basic, and they form the identity matrix. The basis matrix is  $\mathbf{B} = [\mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5] = \mathbf{I}_3$ .

The value of the objective function corresponding to this basic feasible solution is  $z = 0$ . We next compute the reduced cost coefficients corresponding to the nonbasic variables  $x_1$  and  $x_2$ . They are

$$\begin{aligned} r_1 &= c_1 - z_1 = c_1 - (c_3 y_{11} + c_4 y_{21} + c_5 y_{31}) = -2, \\ r_2 &= c_2 - z_2 = c_2 - (c_3 y_{12} + c_4 y_{22} + c_5 y_{32}) = -5. \end{aligned}$$

We would like now to move to an adjacent basic feasible solution for which the objective function value is lower. Naturally, if there is more than one such solution, it is desirable to move to the adjacent basic feasible solution with the lowest objective value. A common practice is to select the most negative value of  $r_j$  and then to bring the corresponding column into the basis (see Exercise 16.18 for an alternative rule for choosing the column to bring into the basis). In our example, we bring  $a_2$  into the basis; that is, we choose  $a_2$  as the new basic column. We then compute  $p = \arg \min\{y_{i0}/y_{i2} : y_{i2} > 0\} = 2$ . We now update the canonical augmented matrix by pivoting about the (2, 2)th entry using the pivot equations:

$$\begin{aligned} y'_{ij} &= y_{ij} - \frac{y_{2j}}{y_{22}} y_{i2}, \quad i \neq 2, \\ y'_{2j} &= \frac{y_{2j}}{y_{22}}. \end{aligned}$$

The resulting updated canonical augmented matrix is

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ \hline 1 & 0 & 1 & 0 & 0 & 4 \\ 0 & 1 & 0 & 1 & 0 & 6 \\ 1 & 0 & 0 & -1 & 1 & 2 \end{array}$$

Note that  $a_2$  entered the basis and  $a_4$  left the basis. The corresponding basic feasible solution is  $\mathbf{x} = [0, 6, 4, 0, 2]^\top$ . We now compute the reduced cost coefficients for the nonbasic columns:

$$\begin{aligned} r_1 &= c_1 - z_1 = -2, \\ r_4 &= c_4 - z_4 = 5. \end{aligned}$$

Because  $r_1 = -2 < 0$ , the current solution is not optimal, and a lower objective function value can be obtained by bringing  $a_1$  into the basis. Proceeding to update the canonical augmented matrix by pivoting about the (3, 1)th element, we obtain

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ \hline 0 & 0 & 1 & 1 & -1 & 2 \\ 0 & 1 & 0 & 1 & 0 & 6 \\ 1 & 0 & 0 & -1 & 1 & 2 \end{array}$$

The corresponding basic feasible solution is  $\mathbf{x} = [2, 6, 2, 0, 0]^\top$ . The reduced cost coefficients are

$$\begin{aligned} r_4 &= c_4 - z_4 = 3, \\ r_5 &= c_5 - z_5 = 2. \end{aligned}$$

Because no reduced cost coefficient is negative, the current basic feasible solution  $\mathbf{x} = [2, 6, 2, 0, 0]^\top$  is optimal. The solution to the original problem is therefore  $x_1 = 2$ ,  $x_2 = 6$ , and the objective function value is 34. ■

We can see from Example 16.2 that we can solve a linear programming problem of any size using the simplex algorithm. To make the calculations in the algorithm more efficient, we discuss the matrix form of the simplex method in the next section.

## 16.5 Matrix Form of the Simplex Method

Consider a linear programming problem in standard form:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Let the first  $m$  columns of  $\mathbf{A}$  be the basic columns. The columns form a square  $m \times m$  nonsingular matrix  $\mathbf{B}$ . The nonbasic columns of  $\mathbf{A}$  form an  $m \times (n-m)$  matrix  $\mathbf{D}$ . We partition the cost vector correspondingly as  $\mathbf{c}^\top = [\mathbf{c}_B^\top, \mathbf{c}_D^\top]$ . Then, the original linear program can be represented as follows:

$$\begin{aligned} & \text{minimize} && \mathbf{c}_B^\top \mathbf{x}_B + \mathbf{c}_D^\top \mathbf{x}_D \\ & \text{subject to} && [\mathbf{B}, \mathbf{D}] \begin{bmatrix} \mathbf{x}_B \\ \mathbf{x}_D \end{bmatrix} = \mathbf{B}\mathbf{x}_B + \mathbf{D}\mathbf{x}_D = \mathbf{b} \\ & && \mathbf{x}_B \geq \mathbf{0}, \mathbf{x}_D \geq \mathbf{0}. \end{aligned}$$

If  $\mathbf{x}_D = \mathbf{0}$ , then the solution  $\mathbf{x} = [\mathbf{x}_B^\top, \mathbf{x}_D^\top]^\top = [\mathbf{x}_B^\top, \mathbf{0}^\top]^\top$  is the basic feasible solution corresponding to the basis  $\mathbf{B}$ . It is clear that for this to be a solution, we need  $\mathbf{x}_B = \mathbf{B}^{-1}\mathbf{b}$ ; that is, the basic feasible solution is

$$\mathbf{x} = \begin{bmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{bmatrix}.$$

The corresponding objective function value is

$$z_0 = \mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{b}.$$

If, on the other hand,  $\mathbf{x}_D \neq \mathbf{0}$ , then the solution  $\mathbf{x} = [\mathbf{x}_B^\top, \mathbf{x}_D^\top]^\top$  is not basic. In this case  $\mathbf{x}_B$  is given by

$$\mathbf{x}_B = \mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x}_D,$$

and the corresponding objective function value is

$$\begin{aligned} z &= \mathbf{c}_B^\top \mathbf{x}_B + \mathbf{c}_D^\top \mathbf{x}_D \\ &= \mathbf{c}_B^\top (\mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{D}\mathbf{x}_D) + \mathbf{c}_D^\top \mathbf{x}_D \\ &= \mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{b} + (\mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{D})\mathbf{x}_D. \end{aligned}$$

Defining

$$\mathbf{r}_D^\top = \mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D},$$

we obtain

$$z = z_0 + \mathbf{r}_D^\top \mathbf{x}_D.$$

The elements of the vector  $\mathbf{r}_D$  are the reduced cost coefficients corresponding to the nonbasic variables.

If  $\mathbf{r}_D \geq \mathbf{0}$ , then the basic feasible solution corresponding to the basis  $\mathbf{B}$  is optimal. If, on the other hand, a component of  $\mathbf{r}_D$  is negative, then the value of the objective function can be reduced by increasing a corresponding component of  $\mathbf{x}_D$ , that is, by changing the basis.

We now use the foregoing observations to develop a matrix form of the simplex method. To this end we first add the cost coefficient vector  $\mathbf{c}^\top$  to the bottom of the augmented matrix  $[\mathbf{A}, \mathbf{b}]$  as follows:

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{c}^\top & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{B} & \mathbf{D} & \mathbf{b} \\ \mathbf{c}_B^\top & \mathbf{c}_D^\top & 0 \end{bmatrix}.$$

We refer to this matrix as the *tableau* of the given LP problem. The tableau contains all relevant information about the linear program.

Suppose that we now apply elementary row operations to the tableau such that the top part of the tableau corresponding to the augmented matrix  $[\mathbf{A}, \mathbf{b}]$  is transformed into canonical form. This corresponds to premultiplying the tableau by the matrix

$$\begin{bmatrix} \mathbf{B}^{-1} & 0 \\ \mathbf{0}^\top & 1 \end{bmatrix}.$$

The result of this operation is

$$\begin{bmatrix} \mathbf{B}^{-1} & 0 \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{B} & \mathbf{D} & \mathbf{b} \\ \mathbf{c}_B^\top & \mathbf{c}_D^\top & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{I}_m & \mathbf{B}^{-1}\mathbf{D} & \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{c}_B^\top & \mathbf{c}_D^\top & 0 \end{bmatrix}.$$

We now apply elementary row operations to the tableau above so that the entries of the last row corresponding to the basic columns become zero. Specifically, this corresponds to premultiplication of the tableau by the matrix

$$\begin{bmatrix} \mathbf{I}_m & 0 \\ -\mathbf{c}_B^\top & 1 \end{bmatrix}.$$

The result is

$$\begin{bmatrix} \mathbf{I}_m & 0 \\ -\mathbf{c}_B^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I}_m & \mathbf{B}^{-1}\mathbf{D} & \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{c}_B^\top & \mathbf{c}_D^\top & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{I}_m & \mathbf{B}^{-1}\mathbf{D} & \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^\top & \mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D} & -\mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{b} \end{bmatrix}.$$

We refer to the resulting tableau as the *canonical tableau corresponding to the basis  $\mathbf{B}$* . Note that the first  $m$  entries of the last column of the canonical

tableau,  $\mathbf{B}^{-1}\mathbf{b}$ , are the values of the basic variables corresponding to the basis  $\mathbf{B}$ . The entries of  $\mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D}$  in the last row are the reduced cost coefficients. The last element in the last row of the tableau,  $-\mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{b}$ , is the negative of the value of the objective function corresponding to the basic feasible solution.

Given an LP problem, we can in general construct many different canonical tableaus, depending on which columns are basic. Suppose that we have a canonical tableau corresponding to a particular basis. Consider the task of computing the tableau corresponding to another basis that differs from the previous basis by a single vector. This can be accomplished by applying elementary row operations to the tableau in a similar fashion as discussed above. We refer to this operation as *updating* the canonical tableau. Note that updating of the tableau involves using exactly the same update equations as we used before in updating the canonical augmented matrix, namely, for  $i = 1, \dots, m+1$ ,

$$\begin{aligned} y'_{ij} &= y_{ij} - \frac{y_{pj}}{y_{pq}} y_{iq}, \quad i \neq p, \\ y'_{pj} &= \frac{y_{pj}}{y_{pq}}, \end{aligned}$$

where  $y_{ij}$  and  $y'_{ij}$  are the  $(i, j)$ th entries of the original and updated canonical tableaus, respectively.

Working with the tableau is a convenient way of implementing the simplex algorithm, since updating the tableau immediately gives us the values of both the basic variables and the reduced cost coefficients. In addition, the (negative of the) value of the objective function can be found in the lower right-hand corner of the tableau. We illustrate the use of the tableau in the following example.

**Example 16.3** Consider the following linear programming problem:

$$\begin{aligned} &\text{maximize} && 7x_1 + 6x_2 \\ &\text{subject to} && 2x_1 + x_2 \leq 3 \\ & && x_1 + 4x_2 \leq 4 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

We first transform the problem into standard form so that the simplex method can be applied. To do this we change the maximization to minimization by multiplying the objective function by  $-1$ . We then introduce two nonnegative slack variables,  $x_3$  and  $x_4$ , and construct the tableau for the problem:

$$\begin{array}{ccccc} & \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ & 2 & 1 & 1 & 0 & 3 \\ & 1 & 4 & 0 & 1 & 4 \\ \mathbf{c}^\top & -7 & -6 & 0 & 0 & 0 \end{array}$$

Notice that this tableau is already in canonical form with respect to the basis  $[a_3, a_4]$ . Hence, the last row contains the reduced cost coefficients, and the rightmost column contains the values of the basic variables. Because  $r_1 = -7$  is the most negative reduced cost coefficient, we bring  $a_1$  into the basis. We then compute the ratios  $y_{10}/y_{11} = 3/2$  and  $y_{20}/y_{21} = 4$ . Because  $y_{10}/y_{11} < y_{20}/y_{21}$ , we get  $p = \arg \min_i \{y_{i0}/y_{i1} : y_{i1} > 0\} = 1$ . We pivot about the  $(1, 1)$ th element of the tableau to obtain

$$\begin{array}{ccccc} 1 & \frac{1}{2} & \frac{1}{2} & 0 & \frac{3}{2} \\ 0 & \frac{7}{2} & -\frac{1}{2} & 1 & \frac{5}{2} \\ 0 & -\frac{5}{2} & \frac{7}{2} & 0 & \frac{21}{2} \end{array}$$

In the second tableau above, only  $r_2$  is negative. Therefore,  $q = 2$  (i.e., we bring  $a_2$  into the basis). Because

$$\frac{y_{10}}{y_{12}} = 3, \quad \frac{y_{20}}{y_{22}} = \frac{5}{7}$$

we have  $p = 2$ . We thus pivot about the  $(2, 2)$ th element of the second tableau to obtain the third tableau:

$$\begin{array}{ccccc} 1 & 0 & \frac{4}{7} & -\frac{1}{7} & \frac{8}{7} \\ 0 & 1 & -\frac{1}{7} & \frac{2}{7} & \frac{5}{7} \\ 0 & 0 & \frac{22}{7} & \frac{5}{7} & \frac{86}{7} \end{array}$$

Because the last row of the third tableau above has no negative elements, we conclude that the basic feasible solution corresponding to the third tableau is optimal. Thus,  $x_1 = 8/7$ ,  $x_2 = 5/7$ ,  $x_3 = 0$ ,  $x_4 = 0$  is the solution to our LP in standard form, and the corresponding objective value is  $-86/7$ . The solution to the original problem is simply  $x_1 = 8/7$ ,  $x_2 = 5/7$ , and the corresponding objective value is  $86/7$ . ■

Degenerate basic feasible solutions may arise in the course of applying the simplex algorithm. In such a situation, the minimum ratio  $y_{i0}/y_{iq}$  is 0. Therefore, even though the basis changes after we pivot about the  $(p, q)$ th element, the basic feasible solution does not (and remains degenerate). It is possible that if we start with a basis corresponding to a degenerate solution, several iterations of the simplex algorithm will involve the same degenerate solution, and eventually the original basis will occur. The entire process will then repeat indefinitely, leading to what is called *cycling*. Such a scenario, although rare in practice, is clearly undesirable. Fortunately, there is a simple rule for choosing  $q$  and  $p$ , due to Bland, that eliminates the cycling problem (see Exercise 16.18):

$$\begin{aligned} q &= \min\{i : r_i < 0\}, \\ p &= \min\{j : y_{j0}/y_{jq} = \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}\}. \end{aligned}$$

## 16.6 Two-Phase Simplex Method

The simplex method requires starting with a tableau for the problem in canonical form; that is, we need an initial basic feasible solution. A brute-force approach to finding a starting basic feasible solution is to choose  $m$  basic columns arbitrarily and transform the tableau for the problem into canonical form. If the rightmost column is positive, then we have a legitimate (initial) basic feasible solution. Otherwise, we would have to pick another candidate basis. Potentially, this brute-force procedure requires  $\binom{n}{m}$  tries, and is therefore not practical.

Certain LP problems have obvious initial basic feasible solutions. For example, if we have constraints of the form  $\mathbf{A}\mathbf{x} \leq \mathbf{b}$  and we add  $m$  slack variables  $z_1, \dots, z_m$ , then the constraints in standard form become

$$[\mathbf{A}, \mathbf{I}_m] \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} = \mathbf{b}, \quad \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \geq \mathbf{0},$$

where  $\mathbf{z} = [z_1, \dots, z_m]^\top$ . The obvious initial basic feasible solution is

$$\begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix},$$

and the basic variables are the slack variables. This was the case in the example in Section 16.5.

Suppose that we are given a linear program in standard form:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

In general, an initial basic feasible solution is not always apparent. We therefore need a systematic method for finding an initial basic feasible solution for general LP problems so that the simplex method can be initialized. For this purpose, suppose that we are given an LP problem in standard form. Consider the following associated *artificial problem*:

$$\begin{aligned} & \text{minimize} && y_1 + y_2 + \cdots + y_m \\ & \text{subject to} && [\mathbf{A}, \mathbf{I}_m] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{b} \\ & && \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{y} = [y_1, \dots, y_m]^\top$ . We call  $\mathbf{y}$  the vector of *artificial variables*. Note that the artificial problem has an obvious initial basic feasible solution:

$$\begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix}.$$

We can therefore solve this problem by the simplex method.

**Proposition 16.1** *The original LP problem has a basic feasible solution if and only if the associated artificial problem has an optimal feasible solution with objective function value zero.*  $\square$

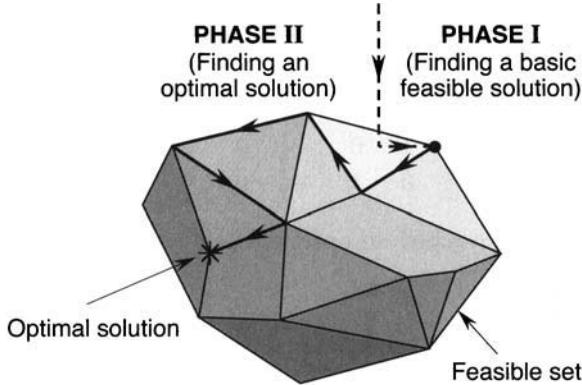
*Proof.*  $\Rightarrow$ : If the original problem has a basic feasible solution  $\mathbf{x}$ , then the vector  $[\mathbf{x}^\top, \mathbf{0}^\top]^\top$  is a basic feasible solution to the artificial problem. Clearly, this solution has an objective function value of zero. This solution is therefore optimal for the artificial problem, since there can be no feasible solution with negative objective function value.

$\Leftarrow$ : Suppose that the artificial problem has an optimal feasible solution with objective function value zero. Then, this solution must have the form  $[\mathbf{x}^\top, \mathbf{0}^\top]^\top$ , where  $\mathbf{x} \geq \mathbf{0}$ . Hence, we have  $\mathbf{Ax} = \mathbf{b}$ , and  $\mathbf{x}$  is a feasible solution to the original problem. By the fundamental theorem of LP, there also exists a basic feasible solution. ■

Assume that the original LP problem has a basic feasible solution. Suppose that the simplex method applied to the associated artificial problem has terminated with an objective function value of zero. Then, as indicated in the proof above, the solution to the artificial problem will have all  $y_i = 0$ ,  $i = 1, \dots, m$ . Hence, assuming nondegeneracy, the basic variables are in the first  $n$  components; that is, none of the artificial variables are basic. Therefore, the first  $n$  components form a basic feasible solution to the original problem. We can then use this basic feasible solution (resulting from the artificial problem) as the initial basic feasible solution for the original LP problem (after deleting the components corresponding to artificial variables). Thus, using artificial variables, we can attack a general linear programming problem by applying the *two-phase simplex method*. In phase I we introduce artificial variables and the artificial objective function and find a basic feasible solution. In phase II we use the basic feasible solution resulting from phase I to initialize the simplex algorithm to solve the original LP problem. The two-phase simplex method is illustrated in Figure 16.1.

**Example 16.4** Consider the following linear programming problem:

$$\begin{aligned} &\text{minimize} && 2x_1 + 3x_2 \\ &\text{subject to} && 4x_1 + 2x_2 \geq 12 \\ & && x_1 + 4x_2 \geq 6 \\ & && x_1, x_2 \geq 0. \end{aligned}$$



**Figure 16.1** Illustration of the two-phase simplex method.

First, we express the problem in standard form by introducing surplus variables:

$$\begin{aligned} & \text{minimize} && 2x_1 + 3x_2 \\ & \text{subject to} && 4x_1 + 2x_2 - x_3 = 12 \\ & && x_1 + 4x_2 - x_4 = 6 \\ & && x_1, \dots, x_4 \geq 0. \end{aligned}$$

For the LP problem above there is no obvious basic feasible solution that we can use to initialize the simplex method. Therefore, we use the two-phase method.

*Phase I.* We introduce artificial variables  $x_5, x_6 \geq 0$ , and an artificial objective function  $x_5 + x_6$ . We form the corresponding tableau for the problem:

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
	4	2	-1	0	1	0	12
	1	4	0	-1	0	1	6
$c^T$	0	0	0	0	1	1	0

To initiate the simplex procedure, we must update the last row of this tableau to transform it into canonical form. We obtain

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
	4	2	-1	0	1	0	12
	1	4	0	-1	0	1	6
	-5	-6	1	1	0	0	-18

The basic feasible solution corresponding to this tableau is not optimal. Therefore, we proceed with the simplex method to obtain the next tableau:

$$\begin{array}{ccccccc} \frac{7}{2} & 0 & -1 & \frac{1}{2} & 1 & -\frac{1}{2} & 9 \\ \frac{1}{4} & 1 & 0 & -\frac{1}{4} & 0 & \frac{1}{4} & \frac{3}{2} \\ -\frac{7}{2} & 0 & 1 & -\frac{1}{2} & 0 & \frac{3}{2} & -9 \end{array}$$

We still have not yet reached an optimal basic feasible solution. Performing another iteration, we get

$$\begin{array}{ccccccc} 1 & 0 & -\frac{2}{7} & \frac{1}{7} & \frac{2}{7} & -\frac{1}{7} & \frac{18}{7} \\ 0 & 1 & \frac{1}{14} & -\frac{2}{7} & -\frac{1}{14} & \frac{2}{7} & \frac{6}{7} \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{array}$$

Both of the artificial variables have been driven out of the basis, and the current basic feasible solution is optimal. We now proceed to phase II.

*Phase II.* We start by deleting the columns corresponding to the artificial variables in the last tableau in phase I and revert back to the original objective function. We obtain

$$\begin{array}{ccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ 1 & 0 & -\frac{2}{7} & \frac{1}{7} & \frac{18}{7} \\ 0 & 1 & \frac{1}{14} & -\frac{2}{7} & \frac{6}{7} \\ \mathbf{c}^\top & 2 & 3 & 0 & 0 \end{array}$$

We transform the last row so that the zeros appear in the basis columns; that is, we transform the tableau above into canonical form:

$$\begin{array}{ccccc} 1 & 0 & -\frac{2}{7} & \frac{1}{7} & \frac{18}{7} \\ 0 & 1 & \frac{1}{14} & -\frac{2}{7} & \frac{6}{7} \\ 0 & 0 & \frac{5}{14} & \frac{4}{7} & -\frac{54}{7} \end{array}$$

All the reduced cost coefficients are nonnegative. Hence, the optimal solution is

$$\mathbf{x} = \left[ \frac{18}{7}, \frac{6}{7}, 0, 0 \right]^\top$$

and the optimal cost is  $54/7$ . ■

## 16.7 Revised Simplex Method

Consider an LP problem in standard form with a matrix  $\mathbf{A}$  of size  $m \times n$ . Suppose that we use the simplex method to solve the problem. Experience suggests that if  $m$  is much smaller than  $n$ , then, in most instances, pivots will

occur in only a small fraction of the columns of the matrix  $\mathbf{A}$ . The operation of pivoting involves updating all the columns of the tableau. However, if a particular column of  $\mathbf{A}$  never enters any basis during the entire simplex procedure, then computations performed on this column are never used. Therefore, if  $m$  is much smaller than  $n$ , the effort expended on performing operations on many of the columns of  $\mathbf{A}$  may be wasted. The *revised simplex method* reduces the amount of computation leading to an optimal solution by eliminating operations on columns of  $\mathbf{A}$  that do not enter the bases.

To be specific, suppose that we are at a particular iteration in the simplex algorithm. Let  $\mathbf{B}$  be the matrix composed of the columns of  $\mathbf{A}$  forming the current basis, and let  $\mathbf{D}$  be the matrix composed of the remaining columns of  $\mathbf{A}$ . The sequence of elementary row operations on the tableau leading to this iteration (represented by matrices  $\mathbf{E}_1, \dots, \mathbf{E}_k$ ) corresponds to premultiplying  $\mathbf{B}, \mathbf{D}$ , and  $\mathbf{b}$  by  $\mathbf{B}^{-1} = \mathbf{E}_k \cdots \mathbf{E}_1$ . In particular, the vector of current values of the basic variables is  $\mathbf{B}^{-1}\mathbf{b}$ . Observe that computation of the current basic feasible solution does not require computation of  $\mathbf{B}^{-1}\mathbf{D}$ ; all we need is the matrix  $\mathbf{B}^{-1}$ . In the revised simplex method we do not compute  $\mathbf{B}^{-1}\mathbf{D}$ . Instead, we only keep track of the basic variables and the revised tableau, which is the tableau  $[\mathbf{B}^{-1}, \mathbf{B}^{-1}\mathbf{b}]$ . Note that this tableau is only of size  $m \times (m+1)$  [compared to the tableau in the original simplex method, which is  $m \times (n+1)$ ]. To see how to update the revised tableau, suppose that we choose the column  $\mathbf{a}_q$  to enter the basis. Let  $\mathbf{y}_q = \mathbf{B}^{-1}\mathbf{a}_q$ ,  $\mathbf{y}_0 = [y_{01}, \dots, y_{0m}]^\top = \mathbf{B}^{-1}\mathbf{b}$ , and  $p = \arg \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\}$  (as in the original simplex method). Then, to update the revised tableau, we form the augmented revised tableau  $[\mathbf{B}^{-1}, \mathbf{y}_0, \mathbf{y}_q]$  and pivot about the  $p$ th element of the last column. We claim that the first  $m+1$  columns of the resulting matrix comprise the updated revised tableau (i.e., we simply remove the last column of the updated augmented revised tableau to obtain the updated revised tableau). To see this, write  $\mathbf{B}^{-1}$  as  $\mathbf{B}^{-1} = \mathbf{E}_k \cdots \mathbf{E}_1$ , and let the matrix  $\mathbf{E}_{k+1}$  represent the pivoting operation above (i.e.,  $\mathbf{E}_{k+1}\mathbf{y}_q = \mathbf{e}_p$ , the  $p$ th column of the  $m \times m$  identity matrix). The matrix  $\mathbf{E}_{k+1}$  is given by

$$\mathbf{E}_{k+1} = \begin{bmatrix} 1 & -y_{1q}/y_{pq} & 0 \\ \ddots & \vdots & \\ & 1/y_{pq} & \\ & \vdots & \ddots \\ 0 & -y_{mq}/y_{pq} & 1 \end{bmatrix}.$$

Then, the updated augmented tableau resulting from the pivoting operation above is  $[\mathbf{E}_{k+1}\mathbf{B}^{-1}, \mathbf{E}_{k+1}\mathbf{y}_0, \mathbf{e}_p]$ . Let  $\mathbf{B}_{\text{new}}$  be the new basis. Then, we have  $\mathbf{B}_{\text{new}}^{-1} = \mathbf{E}_{k+1} \cdots \mathbf{E}_1$ . But notice that  $\mathbf{B}_{\text{new}}^{-1} = \mathbf{E}_{k+1}\mathbf{B}^{-1}$ , and the values of the basic variables corresponding to  $\mathbf{B}_{\text{new}}$  are given by  $\mathbf{y}_{0\text{new}} = \mathbf{E}_{k+1}\mathbf{y}_0$ . Hence, the updated tableau is indeed  $[\mathbf{B}_{\text{new}}^{-1}, \mathbf{y}_{0\text{new}}] = [\mathbf{E}_{k+1}\mathbf{B}^{-1}, \mathbf{E}_{k+1}\mathbf{y}_0]$ .

We summarize the foregoing discussion in the following algorithm.

### Revised Simplex Method

1. Form a revised tableau corresponding to an initial basic feasible solution  $[\mathbf{B}^{-1}, \mathbf{y}_0]$ .
2. Calculate the current reduced cost coefficients vector via

$$\mathbf{r}_D^\top = \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D},$$

where

$$\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1}.$$

3. If  $r_j \geq 0$  for all  $j$ , stop—the current basic feasible solution is optimal.
4. Select a  $q$  such that  $r_q < 0$  (e.g., the  $q$  corresponding to the most negative  $r_q$ ), and compute

$$\mathbf{y}_q = \mathbf{B}^{-1} \mathbf{a}_q.$$

5. If no  $y_{iq} > 0$ , stop—the problem is unbounded; else, compute  $p = \arg \min_i \{y_{io}/y_{iq} : y_{iq} > 0\}$ .
6. Form the augmented revised tableau  $[\mathbf{B}^{-1}, \mathbf{y}_0, \mathbf{y}_q]$ , and pivot about the  $p$ th element of the last column. Form the updated revised tableau by taking the first  $m+1$  columns of the resulting augmented revised tableau (i.e., remove the last column).
7. Go to step 2.

The reason for computing  $\mathbf{r}_D$  in two steps as indicated in step 2 is as follows. We first note that  $\mathbf{r}_D = \mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D}$ . To compute  $\mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D}$ , we can do the multiplication in the order either  $(\mathbf{c}_B^\top \mathbf{B}^{-1}) \mathbf{D}$  or  $\mathbf{c}_B^\top (\mathbf{B}^{-1} \mathbf{D})$ . The former involves two vector-matrix multiplications, whereas the latter involves a matrix-matrix multiplication followed by a vector-matrix multiplication. Clearly, the former is more efficient.

As in the original simplex method, we can use the two-phase method to solve a given LP problem using the revised simplex method. In particular, we use the revised tableau from the final step of phase I as the initial revised tableau in phase II. We illustrate the method in the following example.

**Example 16.5** Consider solving the following LP problem using the revised simplex method:

$$\begin{aligned} &\text{maximize} && 3x_1 + 5x_2 \\ &\text{subject to} && x_1 + x_2 \leq 4 \\ & && 5x_1 + 3x_2 \geq 8 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

First, we express the problem in standard form by introducing one slack and one surplus variable, to obtain

$$\begin{aligned} \text{minimize} \quad & -3x_1 - 5x_2 \\ \text{subject to} \quad & x_1 + x_2 + x_3 = 4 \\ & 5x_1 + 3x_2 - x_4 = 8 \\ & x_1, \dots, x_4 \geq 0. \end{aligned}$$

There is no obvious basic feasible solution to this LP problem. Therefore, we use the two-phase method.

*Phase I.* We introduce one artificial variable  $x_5$  and an artificial objective function  $x_5$ . The tableau for the artificial problem is

$$\begin{array}{ccccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 1 & 1 & 1 & 0 & 0 & 4 \\ 5 & 3 & 0 & -1 & 1 & 8 \\ \mathbf{c}^\top & 0 & 0 & 0 & 0 & 1 & 0 \end{array}$$

We start with an initial basic feasible solution and corresponding  $\mathbf{B}^{-1}$ , as shown in the following revised tableau:

$$\begin{array}{c|ccc} \text{Variable} & \mathbf{B}^{-1} & \mathbf{y}_0 \\ \hline x_3 & 1 & 0 & 4 \\ x_5 & 0 & 1 & 8 \end{array}$$

We compute

$$\begin{aligned} \boldsymbol{\lambda}^\top &= \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, 1], \\ \mathbf{r}_D^\top &= \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [0, 0, 0] - [5, 3, -1] = [-5, -3, 1] = [r_1, r_2, r_4]. \end{aligned}$$

Because  $r_1$  is the most negative reduced cost coefficient, we bring  $\mathbf{a}_1$  into the basis. To do this, we first compute  $\mathbf{y}_1 = \mathbf{B}^{-1} \mathbf{a}_1$ . In this case,  $\mathbf{y}_1 = \mathbf{a}_1$ . We get the augmented revised tableau:

$$\begin{array}{c|ccc|c} \text{Variable} & \mathbf{B}^{-1} & \mathbf{y}_0 & \mathbf{y}_1 \\ \hline x_3 & 1 & 0 & 4 & 1 \\ x_5 & 0 & 1 & 8 & 5 \end{array}$$

We then compute  $p = \arg \min_i \{y_{i0}/y_{iq} : y_{iq} > 0\} = 2$  and pivot about the second element of the last column to get the updated revised tableau:

$$\begin{array}{c|ccc|c} \text{Variable} & \mathbf{B}^{-1} & \mathbf{y}_0 \\ \hline x_3 & 1 & -\frac{1}{5} & \frac{12}{5} \\ x_1 & 0 & \frac{1}{5} & \frac{8}{5} \end{array}$$

We next compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, 0], \\ \mathbf{r}_D^\top &= \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [0, 0, 1] = [r_2, r_4, r_5] \geq \mathbf{0}^\top.\end{aligned}$$

The reduced cost coefficients are all nonnegative. Hence, the solution to the artificial problem is  $[8/5, 0, 12/5, 0, 0]^\top$ . The initial basic feasible solution for phase II is therefore  $[8/5, 0, 12/5, 0]^\top$ .

*Phase II.* The tableau for the original problem (in standard form) is

$$\begin{array}{ccccc|c} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ 1 & 1 & 1 & 0 & 4 \\ 5 & 3 & 0 & -1 & 8 \\ \mathbf{c}^\top & -3 & -5 & 0 & 0 & 0 \end{array}$$

As the initial revised tableau for phase II, we take the final revised tableau from phase I. We then compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, -3] \begin{bmatrix} 1 & -\frac{1}{5} \\ 0 & \frac{1}{5} \end{bmatrix} = \left[0, -\frac{3}{5}\right], \\ \mathbf{r}_D^\top &= \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [-5, 0] - \left[0, -\frac{3}{5}\right] \begin{bmatrix} 1 & 0 \\ 3 & -1 \end{bmatrix} = \left[-\frac{16}{5}, -\frac{3}{5}\right] = [r_2, r_4].\end{aligned}$$

We bring  $\mathbf{a}_2$  into the basis, and compute  $\mathbf{y}_2 = \mathbf{B}^{-1} \mathbf{a}_2$  to get

$$\begin{array}{ccccc} \text{Variable} & \mathbf{B}^{-1} & \mathbf{y}_0 & \mathbf{y}_2 \\ \hline x_3 & 1 & -\frac{1}{5} & \frac{12}{5} & \frac{2}{5} \\ x_1 & 0 & \frac{1}{5} & \frac{8}{5} & \frac{3}{5} \end{array}$$

In this case we get  $p = 2$ . We update this tableau by pivoting about the second element of the last column to get

$$\begin{array}{ccccc} \text{Variable} & \mathbf{B}^{-1} & \mathbf{y}_0 \\ \hline x_3 & 1 & -\frac{1}{3} & \frac{4}{3} \\ x_2 & 0 & \frac{1}{3} & \frac{8}{3} \end{array}$$

We compute

$$\begin{aligned}\boldsymbol{\lambda}^\top &= \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, -5] \begin{bmatrix} 1 & -\frac{1}{3} \\ 0 & \frac{1}{3} \end{bmatrix} = \left[0, -\frac{5}{3}\right], \\ \mathbf{r}_D^\top &= \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [-3, 0] - \left[0, -\frac{5}{3}\right] \begin{bmatrix} 1 & 0 \\ 5 & -1 \end{bmatrix} = \left[\frac{16}{3}, -\frac{5}{3}\right] = [r_1, r_4].\end{aligned}$$

We now bring  $a_4$  into the basis:

Variable	$\mathbf{B}^{-1}$	$\mathbf{y}_0$	$\mathbf{y}_4$
$x_3$	1	$-\frac{1}{3}$	$\frac{4}{3}$
$x_2$	0	$\frac{1}{3}$	$-\frac{1}{3}$

We update the tableau to obtain

Variable	$\mathbf{B}^{-1}$	$\mathbf{y}_0$
$x_4$	3	-1
$x_2$	1	0

We compute

$$\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1} = [0, -5] \begin{bmatrix} 3 & -1 \\ 1 & 0 \end{bmatrix} = [-5, 0],$$

$$\mathbf{r}_D^\top = \mathbf{c}_D^\top - \boldsymbol{\lambda}^\top \mathbf{D} = [-3, 0] - [-5, 0] \begin{bmatrix} 1 & 1 \\ 5 & 0 \end{bmatrix} = [2, 5] = [r_1, r_3].$$

The reduced cost coefficients are all positive. Hence,  $[0, 4, 0, 4]^\top$  is optimal. The optimal solution to the original problem is  $[0, 4]^\top$ . ■

## EXERCISES

**16.1** This question is concerned with elementary row operations and rank.

- a. For the matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -1 & 3 & 2 \\ 2 & -1 & 3 & 0 & 1 \\ 3 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 1 & 1 \end{bmatrix},$$

find its rank by first transforming the matrix using elementary row operations into an upper triangular form.

- b. Find the rank of the following matrix for different values of the parameter  $\gamma$  by first transforming the matrix using elementary row operations into an upper triangular form:

$$\mathbf{A} = \begin{bmatrix} 1 & \gamma & -1 & 2 \\ 2 & -1 & \gamma & 5 \\ 1 & 10 & -6 & 1 \end{bmatrix}.$$

**16.2** Consider the following standard form LP problem:

$$\begin{aligned} & \text{minimize} && 2x_1 - x_2 - x_3 \\ & \text{subject to} && 3x_1 + x_2 + x_4 = 4 \\ & && 6x_1 + 2x_2 + x_3 + x_4 = 5 \\ & && x_1, x_2, x_3, x_4 \geq 0. \end{aligned}$$

- a. Write down the  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  matrices/vectors for the problem.
- b. Consider the basis consisting of the third and fourth columns of  $\mathbf{A}$ , ordered according to  $[a_4, a_3]$ . Compute the canonical tableau corresponding to this basis.
- c. Write down the basic feasible solution corresponding to the basis above, and its objective function value.
- d. Write down the values of the reduced cost coefficients (for all the variables) corresponding to the basis.
- e. Is the basic feasible solution in part c an optimal feasible solution? If yes, explain why. If not, determine which element of the canonical tableau to pivot about so that the new basic feasible solution will have a lower objective function value.
- f. Suppose that we apply the two-phase method to the problem, and at the end of phase I, the tableau for the artificial problem is

$$\begin{array}{ccccccc} 0 & 0 & -1 & 1 & 2 & -1 & 3 \\ 1 & \frac{1}{3} & \frac{1}{3} & 0 & -\frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{array}$$

Does the original problem have a basic feasible solution? Explain.

- g. From the final tableau for phase I in part f, find the initial canonical tableau for phase II.

**16.3** Use the simplex method to solve the following linear program:

$$\begin{aligned} & \text{maximize} && x_1 + x_2 + 3x_3 \\ & \text{subject to} && x_1 + x_3 = 1 \\ & && x_2 + x_3 = 2 \\ & && x_1, x_2, x_3 \geq 0. \end{aligned}$$

**16.4** Consider the linear program

$$\begin{aligned} & \text{maximize} && 2x_1 + x_2 \\ & \text{subject to} && 0 \leq x_1 \leq 5 \\ & && 0 \leq x_2 \leq 7 \\ & && x_1 + x_2 \leq 9. \end{aligned}$$

Convert the problem to standard form and solve it using the simplex method.

**16.5** Consider a standard form linear programming problem with

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} ? & ? & 0 & 1 \\ ? & ? & 1 & 0 \end{bmatrix}, & \mathbf{b} &= \begin{bmatrix} 5 \\ 6 \end{bmatrix}, \\ \mathbf{c}^\top &= \begin{bmatrix} 8 & 7 & ? & ? \end{bmatrix}, \end{aligned}$$

where the “?” symbols signify unknowns to be determined. Suppose that the canonical tableau corresponding to some basis is

$$\left[ \begin{array}{ccccc} 0 & 1 & 1 & 2 & ? \\ 1 & 0 & 3 & 4 & ? \\ 0 & 0 & -1 & 1 & ? \end{array} \right].$$

- a. Find all entries of  $\mathbf{A}$ .
- b. Find all entries of  $\mathbf{c}$ .
- c. Find the basic feasible solution corresponding to the canonical tableau above.
- d. Find all entries in the rightmost column of the tableau.

**16.6** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && c_1|x_1| + c_2|x_2| + \cdots + c_n|x_n| \\ & \text{subject to} && \mathbf{Ax} = \mathbf{b}. \end{aligned}$$

We can convert this problem into an equivalent standard form linear programming problem by introducing the new variables

$$x_i = x_i^+ - x_i^- \text{ where } x_i^+ \geq 0, \quad x_i^- \geq 0, \quad i = 1, 2, \dots, n$$

and

$$|x_i| = x_i^+ + x_i^-, \quad i = 1, 2, \dots, n$$

(See also Exercise 15.3.) Then we can apply the simplex method to solve the equivalent problem. Explain, in two or three sentences, why we will always

have that only either  $x_i^+$  or  $x_i^-$  can be positive but never both  $x_i^+$  and  $x_i^-$  can be positive. In other words, we will always have,  $x_i^+ x_i^- = 0$ .

**16.7** Suppose that we are given a linear programming problem in standard form (written in the usual notation) and are told that the vector  $\mathbf{x} = [1, 0, 2, 3, 0]^\top$  is a basic feasible solution with corresponding relative cost coefficient vector  $\mathbf{r} = [0, 1, 0, 0, -1]^\top$  and objective function value 6. We are also told that the vector  $[-2, 0, 0, 0, 4]^\top$  lies in the nullspace of  $\mathbf{A}$ .

- Write down the canonical tableau corresponding to the given basic feasible solution above, filling in as many values of entries as possible (use the symbol \* for entries that cannot be determined from the information given). Clearly indicate the dimensions of the tableau.
- Find a feasible solution with an objective function value that is strictly less than 6.

**16.8** Consider a standard form linear programming problem (with the usual  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ ). Suppose that it has the following canonical tableau:

$$\begin{array}{ccccccc} 0 & 1 & 0 & 1 & -1 & 5 \\ 1 & 2 & 0 & 0 & -2 & 6 \\ 0 & 3 & 1 & 0 & -3 & 7 \\ 0 & 4 & 0 & 0 & -4 & 8 \end{array}$$

- Find the basic feasible solution corresponding to this canonical tableau and the corresponding value of the objective function.
- Find all the reduced cost coefficient values associated with the tableau.
- Does the given linear programming problem have feasible solutions with arbitrarily negative objective function values?
- Suppose that column  $a_2$  enters the basis. Find the canonical tableau for the new basis.
- Find a feasible solution with objective function value equal to  $-100$ .
- Find a basis for the nullspace of  $\mathbf{A}$ .

**16.9** Consider the problem

$$\begin{aligned} \text{maximize } & -x_1 - 2x_2 \\ \text{subject to } & x_1 \geq 0 \\ & x_2 \geq 1. \end{aligned}$$

- a. Convert the problem into a standard form linear programming problem.
- b. Use the two-phase simplex method to compute the solution to this problem and the value of the objective function at the optimal solution of the problem.

**16.10** Consider the linear programming problem

$$\begin{aligned} & \text{minimize} && -x_1 \\ & \text{subject to} && x_1 - x_2 = 1 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

- a. Write down the basic feasible solution for  $x_1$  as a basic variable.
- b. Compute the canonical augmented matrix corresponding to the basis in part a.
- c. If we apply the simplex algorithm to this problem, under what circumstance does it terminate? (In other words, which stopping criterion in the simplex algorithm is satisfied?)
- d. Show that in this problem, the objective function can take arbitrarily negative values over the constraint set.

**16.11** Find the solution and the value of the optimal cost for the following problem using the revised simplex method:

$$\begin{aligned} & \text{minimize} && x_1 + x_2 \\ & \text{subject to} && x_1 + 2x_2 \geq 3 \\ & && 2x_1 + x_2 \geq 3 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

*Hint:* Start with  $x_1$  and  $x_2$  as basic variables.

**16.12** Solve the following linear programs using the revised simplex method:

- a. Maximize  $-4x_1 - 3x_2$  subject to

$$\begin{aligned} & 5x_1 + x_2 \geq 11 \\ & -2x_1 - x_2 \leq -8 \\ & x_1 + 2x_2 \geq 7 \\ & x_1, x_2 \geq 0. \end{aligned}$$

- b. Maximize  $6x_1 + 4x_2 + 7x_3 + 5x_4$  subject to

$$\begin{aligned}x_1 + 2x_2 + x_3 + 2x_4 &\leq 20 \\6x_1 + 5x_2 + 3x_3 + 2x_4 &\leq 100 \\3x_1 + 4x_2 + 9x_3 + 12x_4 &\leq 75 \\x_1, x_2, x_3, x_4 &\geq 0.\end{aligned}$$

- 16.13** Consider a standard form linear programming problem with

$$\mathbf{A} = \begin{bmatrix} 0 & 2 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 3 & 1 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 7 \\ 8 \\ 9 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 6 \\ c_2 \\ 4 \\ 5 \end{bmatrix}.$$

Suppose that we are told that the reduced cost coefficient vector corresponding to some basis is  $\mathbf{r}^\top = [0, 1, 0, 0]$ .

- a. Find an optimal feasible solution to the problem.
- b. Find  $c_2$ .

- 16.14** Consider the linear programming problem

$$\begin{aligned}\text{minimize} \quad & c_1 x_1 + c_2 x_2 \\ \text{subject to} \quad & 2x_1 + x_2 = 2 \\ & x_1, x_2 \geq 0,\end{aligned}$$

where  $c_1, c_2 \in \mathbb{R}$ . Suppose that the problem has an optimal feasible solution that is not basic.

- a. Find all basic feasible solutions.
- b. Find all possible values of  $c_1$  and  $c_2$ .
- c. At each basic feasible solution, compute the reduced cost coefficients for all nonbasic variables.

- 16.15** Suppose that we apply the simplex method to a given linear programming problem and obtain the following canonical tableau:

$$\begin{array}{ccccc|c} 0 & \beta & 0 & 1 & 4 \\ 1 & \gamma & 0 & 0 & 5 \\ 0 & -3 & 1 & 0 & 6 \\ 0 & 2 - \alpha & 0 & 0 & \delta \end{array}$$

For each of the following conditions, find the set of all parameter values  $\alpha, \beta, \gamma, \delta$  that satisfy the condition.

- The problem has no solution because the objective function values are unbounded.
- The current basic feasible solution is optimal, and the corresponding objective function value is 7.
- The current basic feasible solution is not optimal, and the objective function value strictly decreases if we remove the first column of  $\mathbf{A}$  from the basis.

**16.16** You are given a linear programming problem in standard form. Suppose that you use the two-phase simplex method and arrive at the following canonical tableau in phase I:

$$\left[ \begin{array}{ccccccc|c} ? & 0 & 1 & 1 & ? & ? & 0 & 6 \\ ? & 0 & 0 & ? & ? & ? & 1 & \alpha \\ ? & 1 & 0 & ? & ? & ? & 0 & 5 \\ \gamma & 0 & 0 & \delta & ? & ? & \beta & 0 \end{array} \right].$$

The variables  $\alpha, \beta, \gamma$ , and  $\delta$  are unknowns to be determined. Those entries marked with “?” are unspecified. The only thing you are told is that the value of  $\gamma$  is either 2 or  $-1$ .

- Determine the values of  $\alpha, \beta, \gamma$ , and  $\delta$ .
- Does the given linear programming problem have a feasible solution? If yes, find it. If not, explain why.

**16.17** Suppose we are given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and a vector  $\mathbf{b} \in \mathbb{R}^m$  such that  $\mathbf{b} \geq \mathbf{0}$ . We are interested in an algorithm that, given this  $\mathbf{A}$  and  $\mathbf{b}$ , is guaranteed to produce one of following two outputs: (1) If there exists  $\mathbf{x}$  such that  $\mathbf{Ax} \geq \mathbf{b}$ , then the algorithm produces one such  $\mathbf{x}$ . (2) If no such  $\mathbf{x}$  exists, then the algorithm produces an output to declare so.

Describe in detail how to design this algorithm based on the simplex method.

**16.18** Consider the following linear programming problem (attributed to Beale—see [42, p. 43]):

$$\begin{aligned} \text{minimize} \quad & -\frac{3}{4}x_4 + 20x_5 - \frac{1}{2}x_6 + 6x_7 \\ \text{subject to} \quad & x_1 + \frac{1}{4}x_4 - 8x_5 - x_6 + 9x_7 = 0 \\ & x_2 + \frac{1}{2}x_4 - 12x_5 - \frac{1}{2}x_6 + 3x_7 = 0 \\ & x_3 + x_6 = 1 \\ & x_1, \dots, x_7 \geq 0. \end{aligned}$$

- a. Apply the simplex algorithm to the problem using the rule that  $q$  is the index corresponding to the most negative  $r_q$ . (As usual, if more than one index  $i$  minimizes  $y_{i0}/y_{iq}$ , let  $p$  be the smallest such index.) Start with  $x_1$ ,  $x_2$ , and  $x_3$  as initial basic variables. Notice that cycling occurs.
- b. Repeat part a using *Bland's rule* for choosing  $q$  and  $p$ :

$$\begin{aligned} q &= \min\{i : r_i < 0\}, \\ p &= \min\{j : y_{j0}/y_{jq} = \min_i\{y_{i0}/y_{iq} : y_{iq} > 0\}\}. \end{aligned}$$

Note that Bland's rule for choosing  $p$  corresponds to our usual rule that if more than one index  $i$  minimizes  $y_{i0}/y_{iq}$ , we let  $p$  be the smallest such index.

**16.19** Consider a standard form LP problem. Suppose that we start with an initial basic feasible solution  $\mathbf{x}^{(0)}$  and apply one iteration of the simplex algorithm to obtain  $\mathbf{x}^{(1)}$ .

We can express  $\mathbf{x}^{(1)}$  in terms of  $\mathbf{x}^{(0)}$  as

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)},$$

where  $\alpha_0$  minimizes  $\phi(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$  over all  $\alpha > 0$  such that  $\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}$  is feasible.

- a. Show that  $\mathbf{d}^{(0)} \in \mathcal{N}(\mathbf{A})$ .
- b. As usual, assume that the initial basis is the first  $m$  columns of  $\mathbf{A}$ , and the first iteration involves inserting  $\mathbf{a}_q$  into the basis, where  $q > m$ . Let the  $q$ th column of the canonical augmented matrix be  $\mathbf{y}_q = [y_{1q}, \dots, y_{mq}]^\top$ . Express  $\mathbf{d}^{(0)}$  in terms of  $\mathbf{y}_q$ .

**16.20** Write a simple MATLAB function that implements the simplex algorithm. The inputs are  $\mathbf{c}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{v}$ , where  $\mathbf{v}$  is the vector of indices of basic

columns. Assume that the augmented matrix  $[\mathbf{A}, \mathbf{b}]$  is already in canonical form; that is, the  $v_i$ th column of  $\mathbf{A}$  is  $[0, \dots, 1, \dots, 0]^\top$ , where 1 occurs in the  $i$ th position. The function should output the final solution and the vector of indices of basic columns. Test the MATLAB function on the problem in Example 16.2.

**16.21** Write a MATLAB routine that implements the two-phase simplex method. It may be useful to use the MATLAB function of Exercise 16.20. Test the routine on the problem in Example 16.5.

**16.22** Write a simple MATLAB function that implements the revised simplex algorithm. The inputs are  $\mathbf{c}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{v}$ , and  $\mathbf{B}^{-1}$ , where  $\mathbf{v}$  is the vector of indices of basic columns; that is, the  $i$ th column of  $\mathbf{B}$  is the  $v_i$ th column of  $\mathbf{A}$ . The function should output the final solution, the vector of indices of basic columns, and the final  $\mathbf{B}^{-1}$ . Test the MATLAB function on the problem in Example 16.2.

**16.23** Write a MATLAB routine that implements the two-phase revised simplex method. It may be useful to use the MATLAB function of Exercise 16.22. Test the routine on the problem in Example 16.5.



# CHAPTER 17

---

## DUALITY

---

### 17.1 Dual Linear Programs

Associated with every linear programming problem is a corresponding dual linear programming problem. The dual problem is constructed from the cost and constraints of the original, or primal, problem. Being an LP problem, the dual can be solved using the simplex method. However, as we shall see, the solution to the dual can also be obtained from the solution of the primal problem, and vice versa. Solving an LP problem via its dual may be simpler in certain cases, and also often provides further insight into the nature of the problem. In this chapter we study basic properties of duality and provide an interpretive example of duality. Duality can be used to improve the performance of the simplex algorithm (leading to the primal-dual algorithm), as well as to develop nonsimplex algorithms for solving LP problems (such as Khachiyan's algorithm and Karmarkar's algorithm). We do not discuss this aspect of duality further in this chapter. For an in-depth discussion of the primal-dual method, as well as other aspects of duality, see, for example, [88]. For a description of Khachiyan's algorithm and Karmarkar's algorithm, see Chapter 18.

Suppose that we are given a linear programming problem of the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

We refer to the above as the *primal problem*. We define the corresponding *dual problem* as

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

We refer to the variable  $\boldsymbol{\lambda} \in \mathbb{R}^m$  as the *dual vector*. Note that the cost vector  $\mathbf{c}$  in the primal has moved to the constraints in the dual. The vector  $\mathbf{b}$  on the right-hand side of  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  becomes part of the cost in the dual. Thus, the roles of  $\mathbf{b}$  and  $\mathbf{c}$  are reversed. The form of duality defined above is called the *symmetric form of duality*.

To define the dual of an arbitrary linear programming problem, we use the following procedure. First, we convert the given linear programming problem into an equivalent problem of the primal form shown above. Then, using the symmetric form of duality, we construct the dual to the equivalent problem. We call the resulting problem the dual of the original problem.

Note that based on the definition of duality above, the dual of the dual problem is the primal problem. To see this, we first represent the dual problem in the form

$$\begin{aligned} & \text{minimize} && \boldsymbol{\lambda}^\top (-\mathbf{b}) \\ & \text{subject to} && \boldsymbol{\lambda}^\top (-\mathbf{A}) \geq -\mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Therefore, by the symmetric form of duality, the dual to the above is

$$\begin{aligned} & \text{maximize} && (-\mathbf{c}^\top) \mathbf{x} \\ & \text{subject to} && (-\mathbf{A}) \mathbf{x} \leq -\mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Upon rewriting, we get the original primal problem.

Now consider an LP problem in standard form. This form has equality constraints  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . To formulate the corresponding dual problem, we first convert the equality constraints into equivalent inequality constraints. Specifically, observe that  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is equivalent to

$$\begin{aligned} & \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & -\mathbf{A}\mathbf{x} \geq -\mathbf{b}. \end{aligned}$$

Thus, the original problem with the equality constraints can be written in the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \begin{bmatrix} \mathbf{A} \\ -\mathbf{A} \end{bmatrix} \mathbf{x} \geq \begin{bmatrix} \mathbf{b} \\ -\mathbf{b} \end{bmatrix} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

The LP problem above is in the form of the primal problem in the symmetric form of duality. The corresponding dual is therefore

$$\begin{aligned} & \text{maximize} && [\mathbf{u}^\top \mathbf{v}^\top] \begin{bmatrix} \mathbf{b} \\ -\mathbf{b} \end{bmatrix} \\ & \text{subject to} && [\mathbf{u}^\top \mathbf{v}^\top] \begin{bmatrix} \mathbf{A} \\ -\mathbf{A} \end{bmatrix} \leq \mathbf{c}^\top \\ & && \mathbf{u}, \mathbf{v} \geq \mathbf{0}. \end{aligned}$$

After simple manipulation the dual above can be represented as

$$\begin{aligned} & \text{maximize} && (\mathbf{u} - \mathbf{v})^\top \mathbf{b} \\ & \text{subject to} && (\mathbf{u} - \mathbf{v})^\top \mathbf{A} \leq \mathbf{c}^\top \\ & && \mathbf{u}, \mathbf{v} \geq \mathbf{0}. \end{aligned}$$

Let  $\boldsymbol{\lambda} = \mathbf{u} - \mathbf{v}$ . Then, the dual problem becomes

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top. \end{aligned}$$

Note that since  $\boldsymbol{\lambda} = \mathbf{u} - \mathbf{v}$  and  $\mathbf{u}, \mathbf{v} \geq \mathbf{0}$ , the dual vector  $\boldsymbol{\lambda}$  is not restricted to be nonnegative. We have now derived the dual for a primal in standard form. The form of duality above is referred to as the *asymmetric form of duality*.

We summarize the forms of duality in Tables 17.1 and 17.2. Note that in the asymmetric form of duality, the dual of the dual is also the primal. We can show this by reversing the arguments we used to arrive at the asymmetric form of duality and using the symmetric form of duality.

Recall that at the beginning of this chapter we defined the dual of an arbitrary linear programming problem by first transforming the problem into an equivalent problem of the form of the primal in the symmetric form of duality. We then derived the asymmetric form of duality based on the symmetric form. In both forms of duality the dual of the dual is the primal. Therefore, we now have four forms of primal-dual linear programming pairs: Each of the four linear programming problems in Tables 17.1 and 17.2 is a primal in these four pairs. So, given an arbitrary linear programming problem, we can obtain its

dual by converting the problem into any of the four problems in Tables 17.1 and 17.2.

**Example 17.1** Suppose that we are given the given linear programming problem

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}. \end{aligned}$$

This problem is already close to the form of the dual in Table 17.2. In particular, let us rewrite the above as

$$\begin{aligned} & \text{maximize } \mathbf{x}^\top (-\mathbf{c}) \\ & \text{subject to } \mathbf{x}^\top \mathbf{A}^\top \leq \mathbf{b}^\top. \end{aligned}$$

Its associated dual is then given by the primal in Table 17.2, which has the form

$$\begin{aligned} & \text{minimize } \mathbf{b}^\top \boldsymbol{\lambda} \\ & \text{subject to } \mathbf{A}^\top \boldsymbol{\lambda} = -\mathbf{c} \\ & \boldsymbol{\lambda} \geq \mathbf{0}, \end{aligned}$$

which can be written in the equivalent form

$$\begin{aligned} & \text{maximize } -\boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to } \boldsymbol{\lambda}^\top \mathbf{A} = -\mathbf{c}^\top \\ & \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

If we change the sign of the dual variable, we can rewrite the above in a more “natural” form:

$$\begin{aligned} & \text{maximize } \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to } \boldsymbol{\lambda}^\top \mathbf{A} = \mathbf{c}^\top \\ & \boldsymbol{\lambda} \leq \mathbf{0}. \end{aligned}$$

■

**Example 17.2** This example is adapted from [88]. Recall the diet problem (see Example 15.2). We have  $n$  different types of food. Our goal is to create the most economical diet and at the same time meet or exceed nutritional requirements. Specifically, let  $a_{ij}$  be the amount of the  $i$ th nutrient per unit of the  $j$ th food,  $b_i$  the amount of the  $i$ th nutrient required,  $1 \leq i \leq m$ ,  $c_j$  the

**Table 17.1** Symmetric Form of Duality

Primal		Dual	
minimize	$c^\top x$	maximize	$\lambda^\top b$
subject to	$Ax \geq b$	subject to	$\lambda^\top A \leq c^\top$
	$x \geq 0$		$\lambda \geq 0$

**Table 17.2** Asymmetric Form of Duality

Primal		Dual	
minimize	$c^\top x$	maximize	$\lambda^\top b$
subject to	$Ax = b$	subject to	$\lambda^\top A \leq c^\top$
	$x \geq 0$		

cost per unit of the  $j$ th food, and  $x_i$  the number of units of food  $i$  in the diet. Then, the diet problem can be stated as follows:

$$\begin{aligned}
 & \text{minimize} && c_1x_1 + c_2x_2 + \cdots + c_nx_n \\
 & \text{subject to} && a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \geq b_1 \\
 & && a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \geq b_2 \\
 & && \vdots \\
 & && a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \geq b_m \\
 & && x_1, \dots, x_n \geq 0.
 \end{aligned}$$

Now, consider a health food store that sells nutrient pills (all  $m$  types of nutrients are available). Let  $\lambda_i$  be the price of a unit of the  $i$ th nutrient in the form of nutrient pills. Suppose that we purchase nutrient pills from the health food store at this price such that we exactly meet our nutritional requirements. Then,  $\lambda^\top b$  is the total revenue to the store. Note that since prices are nonnegative, we have  $\lambda \geq 0$ . Consider now the task of substituting nutrient pills for natural food. The cost of buying pills to create the nutritional equivalent of the  $i$ th food synthetically is simply  $\lambda_1a_{1i} + \cdots + \lambda_ma_{mi}$ . Because  $c_i$  is the cost per unit of the  $i$ th food, if

$$\lambda_1a_{1i} + \cdots + \lambda_ma_{mi} \leq c_i,$$

then the cost of the unit of the  $i$ th food made synthetically from nutrient pills is less than or equal to the market price of a unit of the real food. Therefore,

for the health food store to be competitive, the following must hold:

$$\lambda_1 a_{11} + \cdots + \lambda_m a_{m1} \leq c_1$$

⋮

$$\lambda_1 a_{1n} + \cdots + \lambda_m a_{mn} \leq c_n.$$

The problem facing the health food store is to choose the prices  $\lambda_1, \dots, \lambda_m$  such that its revenue is maximized. This problem can be stated as

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Note that this is simply the dual of the diet problem. ■

**Example 17.3** Consider the following linear programming problem:

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 + x_3 \\ & \text{subject to} && 2x_1 - x_2 + 7x_3 \leq 6 \\ & && x_1 + 3x_2 + 4x_3 \leq 9 \\ & && 3x_1 + 6x_2 + x_3 \leq 3 \\ & && x_1, x_2, x_3 \geq 0. \end{aligned}$$

Find the corresponding dual problem and solve it.

We first write the primal problem in standard form by introducing slack variables  $x_4, x_5, x_6$ . This primal problem in standard form is

$$\begin{aligned} & \text{minimize} && [\mathbf{c}^\top, \mathbf{0}^\top] \mathbf{x} \\ & \text{subject to} && [\mathbf{A}, \mathbf{I}] \mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{x} = [x_1, \dots, x_6]^\top$  and

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 7 \\ 1 & 3 & 4 \\ 3 & 6 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 6 \\ 9 \\ 3 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} -2 \\ -5 \\ -1 \end{bmatrix}.$$

The corresponding dual problem (asymmetric form) is

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top [\mathbf{A}, \mathbf{I}] \leq [\mathbf{c}^\top, \mathbf{0}^\top]. \end{aligned}$$

Note that the constraints in the dual can be written as

$$\begin{aligned}\boldsymbol{\lambda}^\top \mathbf{A} &\leq \mathbf{c}^\top \\ \boldsymbol{\lambda} &\leq \mathbf{0}.\end{aligned}$$

To solve the dual problem above, we use the simplex method. For this, we need to express the problem in standard form. We substitute  $\boldsymbol{\lambda}$  by  $-\boldsymbol{\lambda}$  and introduce surplus variables to get

$$\begin{aligned}\text{minimize} \quad & 6\lambda_1 + 9\lambda_2 + 3\lambda_3 \\ \text{subject to} \quad & 2\lambda_1 + \lambda_2 + 3\lambda_3 - \lambda_4 = 2 \\ & -\lambda_1 + 3\lambda_2 + 6\lambda_3 - \lambda_5 = 5 \\ & 7\lambda_1 + 4\lambda_2 + \lambda_3 - \lambda_6 = 1 \\ & \lambda_1, \dots, \lambda_6 \geq 0.\end{aligned}$$

There is no obvious basic feasible solution. Thus, we use the two-phase simplex method to solve the problem.

*Phase I.* We introduce artificial variables  $\lambda_7, \lambda_8, \lambda_9$  and the artificial objective function  $\lambda_7 + \lambda_8 + \lambda_9$ . The tableau for the artificial problem is

	$\lambda_1$	$\lambda_2$	$\lambda_3$	$\lambda_4$	$\lambda_5$	$\lambda_6$	$\lambda_7$	$\lambda_8$	$\lambda_9$	$\mathbf{c}$
	2	1	3	-1	0	0	1	0	0	2
	-1	3	6	0	-1	0	0	1	0	5
	7	4	1	0	0	-1	0	0	1	1
Cost	0	0	0	0	0	0	1	1	1	0

We start with an initial feasible solution and corresponding  $\mathbf{B}^{-1}$ :

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$\lambda_7$	1	0	0	2
$\lambda_8$	0	1	0	5
$\lambda_9$	0	0	1	1

We compute

$$\begin{aligned}\mathbf{r}_D^\top &= [0, 0, 0, 0, 0, 0] - [8, 8, 10, -1, -1, -1] = [-8, -8, -10, 1, 1, 1] \\ &= [r_1, r_2, r_3, r_4, r_5, r_6].\end{aligned}$$

Because  $r_3$  is the most negative reduced cost coefficient, we bring the third column into the basis. In this case,  $\mathbf{y}_3 = [3, 6, 1]^\top$ . We have

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_3$
$\lambda_7$	1	0	0	2	3
$\lambda_8$	0	1	0	5	6
$\lambda_9$	0	0	1	1	1

By inspection,  $p = 1$ , so we pivot about the first element of the last column. The updated tableau is

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$\lambda_3$	$\frac{1}{3}$	0	0	$\frac{2}{3}$
$\lambda_8$	-2	1	0	1
$\lambda_9$	$-\frac{1}{3}$	0	1	$\frac{1}{3}$

We compute

$$\mathbf{r}_D^\top = \left[ -\frac{4}{3}, -\frac{14}{3}, -\frac{7}{3}, 1, 1, \frac{10}{3} \right] = [r_1, r_2, r_4, r_5, r_6, r_7].$$

We bring the second column into the basis to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_2$
$\lambda_3$	$\frac{1}{3}$	0	0	$\frac{2}{3}$	$\frac{1}{3}$
$\lambda_8$	-2	1	0	1	1
$\lambda_9$	$-\frac{1}{3}$	0	1	$\frac{1}{3}$	$\frac{11}{3}$

We update the tableau to get

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$\lambda_3$	$\frac{4}{11}$	0	$-\frac{1}{11}$	$\frac{7}{11}$
$\lambda_8$	$-\frac{21}{11}$	1	$-\frac{3}{11}$	$\frac{10}{11}$
$\lambda_2$	$-\frac{1}{11}$	0	$\frac{3}{11}$	$\frac{1}{11}$

We compute

$$\mathbf{r}_D^\top = \left[ \frac{74}{11}, -\frac{21}{11}, 1, -\frac{3}{11}, \frac{32}{11}, \frac{14}{11} \right] = [r_1, r_4, r_5, r_6, r_7, r_9].$$

We bring the fourth column into the basis:

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$	$\mathbf{y}_4$
$\lambda_3$	$\frac{4}{11}$	0	$-\frac{1}{11}$	$\frac{7}{11}$	$-\frac{4}{11}$
$\lambda_8$	$-\frac{21}{11}$	1	$-\frac{3}{11}$	$\frac{10}{11}$	$\frac{21}{11}$
$\lambda_2$	$-\frac{1}{11}$	0	$\frac{3}{11}$	$\frac{1}{11}$	$\frac{1}{11}$

The updated tableau becomes

Variable	$\mathbf{B}^{-1}$			$\mathbf{y}_0$
$\lambda_3$	0	$\frac{4}{21}$	$-\frac{3}{21}$	$\frac{17}{21}$
$\lambda_4$	-1	$\frac{11}{21}$	$-\frac{3}{21}$	$\frac{10}{21}$
$\lambda_2$	0	$-\frac{1}{21}$	$\frac{6}{21}$	$\frac{1}{21}$

We compute

$$\mathbf{r}_D^\top = [0, 0, 0, 1, 1, 1] = [r_1, r_5, r_6, r_7, r_8, r_9].$$

Because all the reduced cost coefficients are nonnegative, we terminate phase I.

*Phase II.* We use the last tableau in phase I (where none of the artificial variables are basic) as the initial tableau in phase II. Note that we now revert back to the original cost of the dual problem in standard form. We compute

$$\mathbf{r}_D^\top = \left[ -\frac{62}{7}, \frac{1}{7}, \frac{15}{7} \right] = [r_1, r_5, r_6].$$

We bring the first column into the basis to obtain the augmented revised tableau

Variable		$\mathbf{B}^{-1}$	$\mathbf{y}_0$	$\mathbf{y}_1$
$\lambda_3$	0	$\frac{4}{21}$	$-\frac{3}{21}$	$\frac{17}{21}$
$\lambda_4$	-1	$\frac{11}{21}$	$-\frac{3}{21}$	$\frac{10}{21}$
$\lambda_2$	0	$-\frac{1}{21}$	$\frac{6}{21}$	$\frac{1}{21}$

We update the tableau to get

Variable		$\mathbf{B}^{-1}$	$\mathbf{y}_0$
$\lambda_3$	0	$\frac{7}{43}$	$\frac{1}{43}$
$\lambda_4$	-1	$\frac{19}{43}$	$\frac{15}{43}$
$\lambda_1$	0	$-\frac{1}{43}$	$\frac{6}{43}$

We compute

$$\mathbf{r}_D^\top = \left[ \frac{186}{43}, \frac{15}{43}, \frac{39}{43} \right] = [r_2, r_5, r_6].$$

Because all the reduced cost coefficients are nonnegative, the current basic feasible solution is optimal for the dual in standard form. Thus, an optimal solution to the original dual problem is

$$\boldsymbol{\lambda} = \left[ -\frac{1}{43}, 0, -\frac{36}{43} \right]^\top.$$

■

## 17.2 Properties of Dual Problems

In this section we present some basic results on dual linear programs. We begin with the weak duality lemma.

**Lemma 17.1 Weak Duality Lemma.** *Suppose that  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible solutions to primal and dual LP problems, respectively (either in the symmetric or asymmetric form). Then,  $\mathbf{c}^\top \mathbf{x} \geq \boldsymbol{\lambda}^\top \mathbf{b}$ .* □

*Proof.* We prove this lemma only for the asymmetric form of duality. The proof for the symmetric form involves only a slight modification (see Exercise 17.1).

Because  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible, we have  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{x} \geq \mathbf{0}$ , and  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$ . Postmultiplying both sides of the inequality  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$  by  $\mathbf{x} \geq \mathbf{0}$  yields  $\boldsymbol{\lambda}^\top \mathbf{A}\mathbf{x} \leq \mathbf{c}^\top \mathbf{x}$ . But  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , hence  $\boldsymbol{\lambda}^\top \mathbf{b} \leq \mathbf{c}^\top \mathbf{x}$ . ■

The weak duality lemma states that a feasible solution to either problem yields a bound on the optimal cost of the other problem. The cost in the dual is never above the cost in the primal. In particular, the optimal cost of the dual is less than or equal to the optimal cost of the primal, that is, “maximum  $\leq$  minimum.” Hence, if the cost of one of the problems is unbounded, then the other problem has no feasible solution. In other words, if “minimum =  $-\infty$ ” or “maximum =  $+\infty$ ,” then the feasible set in the other problem must be empty.

**Example 17.4** Consider the problem

$$\begin{aligned} & \text{minimize } x \\ & \text{subject to } x \leq 1, \end{aligned}$$

which is clearly unbounded. By Example 17.1, the dual is

$$\begin{aligned} & \text{maximize } \lambda \\ & \text{subject to } \lambda = 1 \\ & \quad \lambda \leq 0, \end{aligned}$$

which is clearly infeasible. ■

It follows from the weak duality lemma that if we are given feasible primal and dual solutions with equal cost, then these solutions must be optimal in their respective problems.

**Theorem 17.1** Suppose that  $\mathbf{x}_0$  and  $\boldsymbol{\lambda}_0$  are feasible solutions to the primal and dual, respectively (either in symmetric or asymmetric form). If  $\mathbf{c}^\top \mathbf{x}_0 = \boldsymbol{\lambda}_0^\top \mathbf{b}$ , then  $\mathbf{x}_0$  and  $\boldsymbol{\lambda}_0$  are optimal solutions to their respective problems. □

*Proof.* Let  $\mathbf{x}$  be an arbitrary feasible solution to the primal problem. Because  $\boldsymbol{\lambda}_0$  is a feasible solution to the dual, by the weak duality lemma,  $\mathbf{c}^\top \mathbf{x} \geq \boldsymbol{\lambda}_0^\top \mathbf{b}$ . So, if  $\mathbf{c}^\top \mathbf{x}_0 = \boldsymbol{\lambda}_0^\top \mathbf{b}$ , then  $\mathbf{c}^\top \mathbf{x}_0 = \boldsymbol{\lambda}_0^\top \mathbf{b} \leq \mathbf{c}^\top \mathbf{x}$ . Hence,  $\mathbf{x}_0$  is optimal for the primal.

On the other hand, let  $\boldsymbol{\lambda}$  be an arbitrary feasible solution to the dual problem. Because  $\mathbf{x}_0$  is a feasible solution to the primal, by the weak duality lemma,  $\mathbf{c}^\top \mathbf{x}_0 \geq \boldsymbol{\lambda}^\top \mathbf{b}$ . Therefore, if  $\mathbf{c}^\top \mathbf{x}_0 = \boldsymbol{\lambda}_0^\top \mathbf{b}$ , then  $\boldsymbol{\lambda}^\top \mathbf{b} \leq \mathbf{c}^\top \mathbf{x}_0 = \boldsymbol{\lambda}_0^\top \mathbf{b}$ . Hence,  $\boldsymbol{\lambda}_0$  is optimal for the dual. ■

We can interpret Theorem 17.1 as follows. The primal seeks to minimize its cost, and the dual seeks to maximize its cost. Because the weak duality

lemma states that “maximum  $\leq$  minimum,” each problem “seeks to reach the other.” When their costs are equal for a pair of feasible solutions, both solutions are optimal, and we have “maximum = minimum.”

It turns out that the converse of Theorem 17.1 is also true, that is, “maximum = minimum” always holds. In fact, we can prove an even stronger result, known as the duality theorem.

**Theorem 17.2 Duality Theorem.** *If the primal problem (either in symmetric or asymmetric form) has an optimal solution, then so does the dual, and the optimal values of their respective objective functions are equal.*  $\square$

*Proof.* We first prove the result for the asymmetric form of duality. Assume that the primal has an optimal solution. Then, by the fundamental theorem of LP, there exists an optimal basic feasible solution. As is our usual notation, let  $\mathbf{B}$  be the matrix of the corresponding  $m$  basic columns,  $\mathbf{D}$  the matrix of the  $n-m$  nonbasic columns,  $\mathbf{c}_B$  the vector of elements of  $\mathbf{c}$  corresponding to basic variables,  $\mathbf{c}_D$  the vector of elements of  $\mathbf{c}$  corresponding to nonbasic variables, and  $\mathbf{r}_D$  the vector of reduced cost coefficients. Then, by Theorem 16.2,

$$\mathbf{r}_D^\top = \mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D} \geq \mathbf{0}^\top.$$

Hence,

$$\mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D} \leq \mathbf{c}_D^\top.$$

Define

$$\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1}.$$

Then,

$$\mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{D} = \boldsymbol{\lambda}^\top \mathbf{D} \leq \mathbf{c}_D^\top.$$

We claim that  $\boldsymbol{\lambda}$  is a feasible solution to the dual. To see this, assume for convenience (and without loss of generality) that the basic columns are the first  $m$  columns of  $\mathbf{A}$ . Then,

$$\boldsymbol{\lambda}^\top \mathbf{A} = \boldsymbol{\lambda}^\top [\mathbf{B}, \mathbf{D}] = [\mathbf{c}_B^\top, \boldsymbol{\lambda}^\top \mathbf{D}] \leq [\mathbf{c}_B^\top, \mathbf{c}_D^\top] = \mathbf{c}^\top.$$

Hence,  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$  and thus  $\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1}$  is feasible.

We claim that  $\boldsymbol{\lambda}$  is also an optimal feasible solution to the dual. To see this, note that

$$\boldsymbol{\lambda}^\top \mathbf{b} = \mathbf{c}_B^\top \mathbf{B}^{-1} \mathbf{b} = \mathbf{c}_B^\top \mathbf{x}_B.$$

Thus, by Theorem 17.1,  $\boldsymbol{\lambda}$  is optimal.

We now prove the symmetric case. First, we convert the primal problem for the symmetric form into the equivalent standard form by adding surplus

variables:

$$\begin{aligned} \text{minimize} \quad & [\mathbf{c}^\top, \mathbf{0}^\top] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \\ \text{subject to} \quad & [\mathbf{A}, -\mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{b} \\ & \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \geq \mathbf{0}. \end{aligned}$$

Note that  $\mathbf{x}$  is optimal for the original primal problem if and only if  $[\mathbf{x}^\top, (\mathbf{Ax} - \mathbf{b})^\top]^\top$  is optimal for the primal in standard form. The dual to the primal in standard form is equivalent to the dual to the original primal in symmetric form. Therefore, the result above for the asymmetric case applies also to the symmetric case.

This completes the proof. ■

**Example 17.5** Recall Example 17.2, where we formulated the dual of the diet problem. From the duality theorem, the maximum revenue for the health food store is the *same* as the minimum cost of a diet that satisfies all of the nutritional requirements; that is,  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ . ■

Consider a primal-dual pair in asymmetric form. Suppose that we solve the primal problem using the simplex method. The proof of the duality theorem suggests a way of obtaining an optimal solution to the dual by using the last row of the final simplex tableau for the primal. First, we write the tableau for the primal problem:

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{c}^\top & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{B} & \mathbf{D} & \mathbf{b} \\ \mathbf{c}_B^\top & \mathbf{c}_D^\top & 0 \end{bmatrix}.$$

Suppose that the matrix  $\mathbf{B}$  is the basis for an optimal basic feasible solution. Then, the final simplex tableau is

$$\begin{bmatrix} \mathbf{I} & \mathbf{B}^{-1}\mathbf{D} & \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^\top & \mathbf{r}_D^\top & -\mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{b} \end{bmatrix},$$

where  $\mathbf{r}_D^\top = \mathbf{c}_D^\top - \mathbf{c}_B^\top \mathbf{B}^{-1}\mathbf{D}$ . In the proof of the duality theorem we have shown that  $\boldsymbol{\lambda}^\top = \mathbf{c}_B^\top \mathbf{B}^{-1}$  is an optimal solution to the dual. The vector  $\boldsymbol{\lambda}$  can be obtained from the final tableau above. Specifically, if  $\text{rank } \mathbf{D} = m$ , then we can solve for  $\boldsymbol{\lambda}$  using the vector  $\mathbf{r}_D$ , via the equation

$$\boldsymbol{\lambda}^\top \mathbf{D} = \mathbf{c}_D^\top - \mathbf{r}_D^\top.$$

Of course, it may turn out that  $\text{rank } \mathbf{D} < m$ . In this case as we now show, we have additional linear equations that allow us to solve for  $\boldsymbol{\lambda}$ . To this end, recall that  $\boldsymbol{\lambda}^\top \mathbf{B} = \mathbf{c}_B^\top$ . Therefore, if we define  $\mathbf{r}^\top = [\mathbf{0}^\top, \mathbf{r}_D^\top]$ , then combining the equations  $\boldsymbol{\lambda}^\top \mathbf{D} = \mathbf{c}_D^\top - \mathbf{r}_D^\top$  and  $\boldsymbol{\lambda}^\top \mathbf{B} = \mathbf{c}_B^\top$  yields

$$\boldsymbol{\lambda}^\top \mathbf{A} = \mathbf{c}^\top - \mathbf{r}^\top.$$

The vector  $\boldsymbol{\lambda}$  may be easy to obtain from the equation  $\boldsymbol{\lambda}^\top \mathbf{D} = \mathbf{c}_D^\top - \mathbf{r}_D^\top$  if  $\mathbf{D}$  takes certain special forms. In particular, this is the case if  $\mathbf{D}$  has an  $m \times m$  identity matrix embedded in it; that is, by rearranging the positions of the columns of  $\mathbf{D}$ , if necessary,  $\mathbf{D}$  has the form  $\mathbf{D} = [\mathbf{I}_m, \mathbf{G}]$ , where  $\mathbf{G}$  is an  $m \times (n - 2m)$  matrix. In this case we can write the equation  $\boldsymbol{\lambda}^\top \mathbf{D} = \mathbf{c}_D^\top - \mathbf{r}_D^\top$  as

$$[\boldsymbol{\lambda}^\top, \boldsymbol{\lambda}^\top \mathbf{G}] = [\mathbf{c}_I^\top, \mathbf{c}_G^\top] - [\mathbf{r}_I^\top, \mathbf{r}_G^\top].$$

Hence,  $\boldsymbol{\lambda}$  is given by

$$\boldsymbol{\lambda}^\top = \mathbf{c}_I^\top - \mathbf{r}_I^\top.$$

In other words, the solution to the dual is obtained by subtracting the reduced costs coefficients corresponding to the identity matrix in  $\mathbf{D}$  from the corresponding elements in the vector  $\mathbf{c}$  (i.e.,  $\mathbf{c}_I$ ).

For example, if we have a problem where we introduced slack variables, and the basic variables for the optimal basic feasible solution do not include any of the slack variables, then the matrix  $\mathbf{D}$  has an identity matrix embedded in it. In addition, in this case we have  $\mathbf{c}_I = \mathbf{0}$ . Therefore,  $\boldsymbol{\lambda} = -\mathbf{r}_I$  is an optimal solution to the dual.

**Example 17.6** In Example 17.3, the tableau for the primal in standard form is

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
2	-1	7	1	0	0	6
1	3	4	0	1	0	9
3	6	1	0	0	1	3
$\mathbf{c}^\top$	-2	-5	-1	0	0	0

If we now solve the problem using the simplex method, we get the following final simplex tableau:

$\frac{15}{43}$	0	1	$\frac{6}{43}$	0	$\frac{1}{43}$	$\frac{39}{43}$
$-\frac{74}{43}$	0	0	$-\frac{21}{43}$	1	$-\frac{25}{43}$	$\frac{186}{43}$
$\frac{19}{43}$	1	0	$-\frac{1}{43}$	0	$\frac{7}{43}$	$\frac{15}{43}$
$\mathbf{r}^\top$	$\frac{24}{43}$	0	0	$\frac{1}{43}$	0	$\frac{36}{43}$
						$\frac{114}{43}$

We can now find the solution of the dual from the above simplex tableau using the equation  $\boldsymbol{\lambda}^\top \mathbf{D} = \mathbf{c}_D^\top - \mathbf{r}_D^\top$ :

$$[\lambda_1, \lambda_2, \lambda_3] \begin{bmatrix} 2 & 1 & 0 \\ 1 & 0 & 0 \\ 3 & 0 & 1 \end{bmatrix} = [-2, 0, 0] - \left[ \frac{24}{43}, \frac{1}{43}, \frac{36}{43} \right].$$

Solving the above, we get

$$\boldsymbol{\lambda}^\top = \left[ -\frac{1}{43}, 0, -\frac{36}{43} \right],$$

which agrees with our solution in Example 17.3. ■

We now summarize our discussion relating the solutions of the primal and dual problems. If one has unbounded objective function values, then the other has no feasible solution. If one has an optimal feasible solution, then so does the other (and their objective function values are equal). One final case remains: What can we say if one (the primal, say) has no feasible solution? In this case clearly the other (the dual, say) cannot have an optimal solution. However, is it necessarily the case that the dual is unbounded? The answer is no: If one of the problems has no feasible solution, then the other may or may not have a feasible solution. The following example shows that there exists a primal-dual pair of problems for which both have no feasible solution.

**Example 17.7** Consider the primal problem

$$\begin{aligned} & \text{minimize} && [1, -2]\mathbf{x} \\ & \text{subject to} && \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \mathbf{x} \geq \begin{bmatrix} 2 \\ -1 \end{bmatrix} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

The problem has no feasible solution, because the constraints require that  $x_1 - x_2 \geq 2$  and  $x_1 - x_2 \leq 1$ . Based on symmetric duality, the dual is

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \begin{bmatrix} 2 \\ -1 \end{bmatrix} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \leq [1, -2] \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

The dual also has no feasible solution, because the constraints require that  $\lambda_1 - \lambda_2 \leq 1$  and  $\lambda_1 - \lambda_2 \geq 2$ . ■

We end this chapter by presenting the following theorem, which describes an alternative form of the relationship between the optimal solutions to the primal and dual problems.

**Theorem 17.3 Complementary Slackness Condition.** *The feasible solutions  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  to a dual pair of problems (either in symmetric or asymmetric form) are optimal if and only if:*

1.  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = 0.$
2.  $\boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) = 0.$

□

*Proof.* We first prove the result for the asymmetric case. Note that condition 2 holds trivially for this case. Therefore, we consider only condition 1.

⇒: If the two solutions are optimal, then by Theorem 17.2,  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ . Because  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , we also have  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = 0$ .

⇐: If  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = 0$ , then  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{A}\mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ . Therefore, by Theorem 17.1,  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are optimal.

We now prove the result for the symmetric case.

⇒: We first show condition 1. If the two solutions are optimal, then by Theorem 17.2,  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ . Because  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  and  $\boldsymbol{\lambda} \geq \mathbf{0}$ , we have

$$(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = \mathbf{c}^\top \mathbf{x} - \boldsymbol{\lambda}^\top \mathbf{A}\mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b} - \boldsymbol{\lambda}^\top \mathbf{A}\mathbf{x} = \boldsymbol{\lambda}^\top (\mathbf{b} - \mathbf{A}\mathbf{x}) \leq 0.$$

On the other hand, since  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$  and  $\mathbf{x} \geq \mathbf{0}$ , we have  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} \geq 0$ . Hence,  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = 0$ . To show condition 2, note that since  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  and  $\boldsymbol{\lambda} \geq \mathbf{0}$ , we have  $\boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) \geq 0$ . On the other hand, since  $\boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top$  and  $\mathbf{x} \geq \mathbf{0}$ , we have  $\boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) = (\boldsymbol{\lambda}^\top \mathbf{A} - \mathbf{c}^\top)\mathbf{x} \leq 0$ .

⇐: Combining conditions 1 and 2, we get  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{A}\mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ . Hence, by Theorem 17.1,  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are optimal. ■

Note that if  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible solutions for the dual pair of problems, we can write condition 1, that is,  $(\mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A})\mathbf{x} = 0$ , as “ $x_i > 0$  implies that  $\boldsymbol{\lambda}^\top \mathbf{a}_i = c_i$ ,  $i = 1, \dots, n$ ,” that is, for any component of  $\mathbf{x}$  that is positive, the corresponding constraint for the dual must be an equality at  $\boldsymbol{\lambda}$ . Also, observe that the statement “ $x_i > 0$  implies that  $\boldsymbol{\lambda}^\top \mathbf{a}_i = c_i$ ” is equivalent to “ $\boldsymbol{\lambda}^\top \mathbf{a}_i < c_i$  implies that  $x_i = 0$ .” A similar representation can be written for condition 2.

Consider the asymmetric form of duality. Recall that for the case of an optimal basic feasible solution  $\mathbf{x}$ ,  $\mathbf{r}^\top = \mathbf{c}^\top - \boldsymbol{\lambda}^\top \mathbf{A}$  is the vector of reduced cost coefficients. Therefore, in this case, the complementary slackness condition can be written as  $\mathbf{r}^\top \mathbf{x} = 0$ .

**Example 17.8** Suppose that you have 26 dollars and you wish to purchase some gold. You have a choice of four vendors, with prices (in dollars per ounce) of  $1/2$ ,  $1$ ,  $1/7$ , and  $1/4$ , respectively. You wish to spend your entire 26 dollars by purchasing gold from these four vendors, where  $x_i$  is the dollars you spend on vendor  $i$ ,  $i = 1, 2, 3, 4$ .

- a. Formulate the linear programming problem (in standard form) that reflects your desire to obtain the maximum weight in gold.
- b. Write down the dual of the linear programming problem in part a, and find the solution to the dual.

- c. Use the complementary slackness condition together with part b to find the optimal values of  $x_1, \dots, x_4$ .

*Solution:*

- a. The corresponding linear programming problem is

$$\begin{aligned} & \text{minimize} && -(2x_1 + x_2 + 7x_3 + 4x_4) \\ & \text{subject to} && x_1 + x_2 + x_3 + x_4 = 26 \\ & && x_1, x_2, x_3, x_4 \geq 0. \end{aligned}$$

- b. The dual problem is

$$\begin{aligned} & \text{maximize} && 26\lambda \\ & \text{subject to} && \lambda \leq -2 \\ & && \lambda \leq -1 \\ & && \lambda \leq -7 \\ & && \lambda \leq -4. \end{aligned}$$

The solution is clearly  $\lambda = -7$ . (*Note:* It is equally valid to have a dual problem with variable  $\lambda' = -\lambda$ .)

- c. By the complementary slackness condition, we know that if we can find a vector  $\mathbf{x}$  that is feasible in the primal and satisfies  $(-[2, 1, 7, 4] - (-7)[1, 1, 1, 1])\mathbf{x} = 0$ , then this  $\mathbf{x}$  is optimal in the primal (original) problem. We can rewrite the conditions above as

$$[1, 1, 1, 1]\mathbf{x} = 26, \quad \mathbf{x} \geq \mathbf{0}, \quad [5, 6, 0, 3]\mathbf{x} = 0.$$

By  $\mathbf{x} \geq \mathbf{0}$  and  $[5, 6, 0, 3]\mathbf{x} = 0$ , we conclude that  $x_1 = x_2 = x_4 = 0$ , and by  $[1, 1, 1, 1]\mathbf{x} = 26$  we then conclude that  $\mathbf{x} = [0, 0, 26, 0]^T$ . ■

## EXERCISES

**17.1** Prove the weak duality lemma for the symmetric form of duality.

**17.2** Find the dual of the optimization problem in Exercise 15.8.

**17.3** Consider the following linear program:

$$\begin{aligned} & \text{maximize} && 2x_1 + 3x_2 \\ & \text{subject to} && x_1 + 2x_2 \leq 4 \\ & && 2x_1 + x_2 \leq 5 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

- a. Use the simplex method to solve the problem.
- b. Write down the dual of the linear program and solve the dual.

**17.4** Consider the linear program

$$\begin{aligned} & \text{minimize} && 4x_1 + 3x_2 \\ & \text{subject to} && 5x_1 + x_2 \geq 11 \\ & && 2x_1 + x_2 \geq 8 \\ & && x_1 + 2x_2 \geq 7 \\ & && x_1, x_2 \geq 0. \end{aligned}$$

Write down the corresponding dual problem and find the solution to the dual.  
(Compare this problem with the one in Exercise 16.12, part a.)

**17.5** Consider the following primal problem:

$$\begin{aligned} & \text{maximize} && x_1 + 2x_2 \\ & \text{subject to} && -2x_1 + x_2 + x_3 = 2 \\ & && -x_1 + 2x_2 + x_4 = 7 \\ & && x_1 + x_5 = 3 \\ & && x_i \geq 0, \quad i = 1, 2, 3, 4, 5. \end{aligned}$$

- a. Construct the dual problem corresponding to the primal problem above.
- b. It is known that the solution to the primal above is  $\mathbf{x}^* = [3, 5, 3, 0, 0]^\top$ . Find the solution to the dual.

**17.6** Consider the linear programming problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}. \end{aligned}$$

- a. Find the dual to this problem.
- b. Suppose that  $\mathbf{b} = \mathbf{0}$  and there exists a vector  $\mathbf{y} \geq \mathbf{0}$  such that  $\mathbf{y}^\top \mathbf{A} + \mathbf{c}^\top = \mathbf{0}^\top$ . Does this problem have an optimal feasible solution? If yes, find it. If no, explain why not. Give complete explanations.

**17.7** Convert the following optimization problem into a linear programming problem and solve it:

$$\begin{aligned} & \text{maximize} && -|x_1| - |x_2| - |x_3| \\ & \text{subject to} && \begin{bmatrix} 1 & 1 & -1 \\ 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}. \end{aligned}$$

Then construct its dual program and solve it.

*Hint:* Introduce two sets of nonnegative variables:  $x_i^+ \geq 0$ ,  $x_i^- \geq 0$ . Then represent the optimization problem using the variables above. Note that only one  $x_i^+$  and  $x_i^-$  can be nonzero at a time. If  $x_i \geq 0$  then  $x_i^+ = x_i$  and  $x_i^- = 0$ . On the other hand, if  $x_i < 0$  then  $x_i^+ = 0$  and  $x_i = -x_i^-$ . See Exercise 16.6.

### 17.8 Consider the linear program

$$\begin{aligned} &\text{minimize} && x_1 + \cdots + x_n, \quad x_1, \dots, x_n \in \mathbb{R} \\ &\text{subject to} && a_1 x_1 + \cdots + a_n x_n = 1 \\ & && x_1, \dots, x_n \geq 0, \end{aligned}$$

where  $0 < a_1 < a_2 < \cdots < a_n$ .

- a. Write down the dual to the problem and find a solution to the dual in terms of  $a_1, \dots, a_n$ .
- b. State the duality theorem and use it to find a solution to the primal problem above.
- c. Suppose that we apply the simplex algorithm to the primal problem. Show that if we start at a nonoptimal initial basic feasible solution, the algorithm terminates in one step if and only if we use the rule where the next nonbasic column to enter the basis is the one with the most negative reduced cost coefficient.

### 17.9 You are given the following linear programming problem:

$$\begin{aligned} &\text{maximize} && c_1 x_1 + \cdots + c_n x_n \\ &\text{subject to} && x_1 + \cdots + x_n = 1 \\ & && x_1, \dots, x_n \geq 0, \end{aligned}$$

where  $c_1, \dots, c_n \in \mathbb{R}$  are constants.

- a. Write down the dual linear program for the primal problem.
- b. Suppose you know that  $c_4 > c_i$  for all  $i \neq 4$ . Use this information to solve the dual.
- c. Use part b to solve the linear programming problem.

### 17.10 Consider the linear programming problem

$$\begin{aligned} &\text{maximize} && \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{0} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{c} = [1, 1, \dots, 1]^\top$ . Assume that the problem has a solution.

- a. Write down the dual of this problem.
- b. Find the solution to the problem.
- c. What can you say about the constraint set for the problem?

**17.11** Consider a given linear programming problem in standard form (written in the usual notation).

- a. Write down the associated artificial problem for the problem (used in the two-phase method).
- b. Write down the dual to the artificial problem from part a.
- c. Prove that if the original linear programming problem has a feasible solution, then the dual problem in part b has an optimal feasible solution.

**17.12** Consider a pair of primal and dual linear programming problems (either in symmetric or asymmetric form). Identify which of the following situations are possible (depending on the particular primal-dual pair) and which are impossible (regardless of the primal-dual pair). In each case, justify your answer (citing results such as the weak duality lemma and the duality theorem whenever needed).

- a. The primal has a feasible solution, and the dual has no feasible solution.
- b. The primal has an optimal feasible solution, and the dual has no optimal feasible solution.
- c. The primal has a feasible solution but no optimal feasible solution, and the dual has an optimal feasible solution.

**17.13** Consider an LP problem in standard form. Suppose that  $\mathbf{x}$  is a feasible solution to the problem. Show that if there exist  $\boldsymbol{\lambda}$  and  $\boldsymbol{\mu}$  such that

$$\begin{aligned}\mathbf{A}^\top \boldsymbol{\lambda} + \boldsymbol{\mu} &= \mathbf{c} \\ \boldsymbol{\mu}^\top \mathbf{x} &= 0 \\ \boldsymbol{\mu} &\geq 0,\end{aligned}$$

then  $\mathbf{x}$  is an optimal feasible solution to the problem and  $\boldsymbol{\lambda}$  is an optimal feasible solution to the dual. The conditions above, called the *Karush-Kuhn-Tucker optimality conditions for LP*, are discussed in detail in Chapters 21 and 22.

**17.14** Consider the linear program

$$\begin{aligned} & \text{maximize } \mathbf{c}^\top \mathbf{x}, \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Use the symmetric form of duality to derive the dual of this linear program and show that the constraint in the dual involving  $\mathbf{A}$  can be written as an equality constraint.

*Hint:* Write  $\mathbf{x} = \mathbf{u} - \mathbf{v}$ , with  $\mathbf{u}, \mathbf{v} \geq \mathbf{0}$ .

**17.15** Consider the linear program

$$\begin{aligned} & \text{minimize } x_1 + x_2 \\ & \text{subject to } x_1 + 2x_2 \geq 3 \\ & \quad 2x_1 + x_2 \geq 3 \\ & \quad x_1, x_2 \geq 0. \end{aligned}$$

The solution to the problem is  $[1, 1]^\top$  (see Exercise 16.11). Write down the dual to the problem, solve the dual, and verify that the duality theorem holds.

**17.16** Consider the problem

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n \\ & \text{subject to } \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

For this problem we have the following theorem.

**Theorem:** A solution to the foregoing problem exists if and only if  $\mathbf{c} \geq \mathbf{0}$ . Moreover, if a solution exists,  $\mathbf{0}$  is a solution.

Use the duality theorem to prove this theorem (see also Exercise 22.15).

**17.17** Let  $\mathbf{A}$  be a given matrix and  $\mathbf{b}$  a given vector. Show that there exists a vector  $\mathbf{x}$  such that  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  and  $\mathbf{x} \geq \mathbf{0}$  if and only if for any vector  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ , we have  $\mathbf{b}^\top \mathbf{y} \leq 0$ .

**17.18** Let  $\mathbf{A}$  be a given matrix and  $\mathbf{b}$  a given vector. We wish to prove the following result: There exists a vector  $\mathbf{x}$  such that  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $\mathbf{x} \geq \mathbf{0}$  if and only if for any given vector  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$ , we have  $\mathbf{b}^\top \mathbf{y} \leq 0$ . This result is known as *Farkas's transposition theorem*. Our argument is based on duality theory, consisting of the following parts.

a. Consider the primal linear program

$$\begin{aligned} & \text{minimize } \mathbf{0}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \\ & \quad \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

Write down the dual of this problem using the notation  $\mathbf{y}$  for the dual variable.

- b. Show that the feasible set of the dual problem is guaranteed to be nonempty.

*Hint:* Think about an obvious feasible point.

- c. Suppose that for any  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$ , we have  $\mathbf{b}^\top \mathbf{y} \leq 0$ . In this case what can you say about whether or not the dual has an optimal feasible solution?

*Hint:* Think about the obvious feasible point in part b.

- d. Suppose that for any  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$ , we have  $\mathbf{b}^\top \mathbf{y} \leq 0$ . Use parts b and c to show that there exists  $\mathbf{x}$  such that  $\mathbf{Ax} = \mathbf{b}$  and  $\mathbf{x} \geq \mathbf{0}$ . (This proves one direction of Farkas's transposition theorem.)

- e. Suppose that  $\mathbf{x}$  satisfies  $\mathbf{Ax} = \mathbf{b}$  and  $\mathbf{x} \geq \mathbf{0}$ . Let  $\mathbf{y}$  be an arbitrary vector satisfying  $\mathbf{A}^\top \mathbf{y} \leq \mathbf{0}$ . Show that  $\mathbf{b}^\top \mathbf{y} \leq 0$ . (This proves the other direction of Farkas's transposition theorem.)

**17.19** Let  $\mathbf{A}$  be a given matrix and  $\mathbf{b}$  a given vector. Show that there exists a vector  $\mathbf{x}$  such that  $\mathbf{Ax} \leq \mathbf{b}$  if and only if for any given vector  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} = \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ , we have  $\mathbf{b}^\top \mathbf{y} \geq 0$ . This result is known as *Gale's transposition theorem*.

**17.20** Let  $\mathbf{A}$  be a given matrix. Show that there exists a vector  $\mathbf{x}$  such that  $\mathbf{Ax} < \mathbf{0}$  if and only if for any given vector  $\mathbf{y}$  satisfying  $\mathbf{A}^\top \mathbf{y} = \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ , we have  $\mathbf{y} = \mathbf{0}$  (i.e.,  $\mathbf{y} = \mathbf{0}$  is the only vector satisfying  $\mathbf{A}^\top \mathbf{y} = \mathbf{0}$  and  $\mathbf{y} \geq \mathbf{0}$ ). This result is known as *Gordan's transposition theorem*.

**17.21** Let  $\mathbf{P} \in \mathbb{R}^{n \times n}$  be a matrix with the property that each element is in the real interval  $[0, 1]$ , and the sum of the elements of each row is equal to 1; call such a matrix a *stochastic matrix*. Now consider a vector  $\mathbf{x} \geq \mathbf{0}$  such that  $\mathbf{x}^\top \mathbf{e} = 1$ , where  $\mathbf{e} = [1, \dots, 1]^\top$ ; call such a vector  $\mathbf{x}$  a *probability vector*.

We wish to prove the following result: For any stochastic matrix  $\mathbf{P}$ , there exists a probability vector  $\mathbf{x}$  such that  $\mathbf{x}^\top \mathbf{P} = \mathbf{x}^\top$ . Although this is a key result in probability theory (under the topic of *Markov chains*), our argument is based on duality theory (for linear programming), consisting of the following parts.

- a. Consider the primal linear program:

$$\begin{aligned} &\text{maximize} && \mathbf{x}^\top \mathbf{e} \\ &\text{subject to} && \mathbf{x}^\top \mathbf{P} = \mathbf{x}^\top \\ &&& \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Write down the dual of this problem.

- b. Show that the dual is not feasible (i.e., there does not exist a feasible solution to the dual).
- Hint:* Derive a contradiction based on  $\mathbf{P}\mathbf{y} > \mathbf{y}$ ; think about the largest element of  $\mathbf{y}$  (call it  $y_i$ ).
- c. Is the primal feasible? What can you deduce about whether or not the primal is unbounded?
  - d. Use part c to deduce the desired result: that there exists a vector  $\mathbf{x} \geq \mathbf{0}$  such that  $\mathbf{x}^\top \mathbf{P} = \mathbf{x}^\top$  and  $\mathbf{x}^\top \mathbf{e} = 1$ .

**17.22** Suppose that you are presented with a “black box” that implements a function  $\phi$  defined as follows: Given positive integers  $m$  and  $n$ , a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and a vector  $\mathbf{b} \in \mathbb{R}^m$ , the value of  $\phi(m, n, \mathbf{A}, \mathbf{b})$  is a vector  $\mathbf{x} = \phi(m, n, \mathbf{A}, \mathbf{b})$  that satisfies  $\mathbf{Ax} \geq \mathbf{b}$ , if such a vector exists. In other words, the black box solves a linear *feasibility problem*.

Now, given  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{c} \in \mathbb{R}^n$ , consider the linear programming problem

$$\begin{aligned} &\text{minimize} && \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} && \mathbf{Ax} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Express a solution to this problem in terms of the function  $\phi$  given above. In other words, show how we can use the black box to solve this linear programming problem.

*Hint:* Find the appropriate inputs to the black box such that the output immediately gives a solution to the linear programming problem. You should use the black box only once.

**17.23** This exercise illustrates the use of duality to compute the sensitivity of the optimal objective function value with respect to perturbations in the constraint.

Consider a primal linear programming problem and its dual (in either symmetric or asymmetric form). Let us view the  $\mathbf{b}$  vector in the primal as a parameter that we can vary, and that we wish to calculate the change in the optimal objective function value if we perturb  $\mathbf{b}$  by a small perturbation  $\Delta\mathbf{b}$  (i.e., replace  $\mathbf{b}$  by  $\mathbf{b} + \Delta\mathbf{b}$ ).

- a. To make the problem precise, let  $z(\mathbf{b})$  be the optimal value of the primal objective function. Let  $\boldsymbol{\lambda}$  denote the corresponding optimal dual vector. Calculate the gradient of  $z$  at  $\mathbf{b}$ :  $\nabla z(\mathbf{b})$ . Write the answer in terms of  $\boldsymbol{\lambda}$ . You may assume that the optimal dual vector remains fixed in a neighborhood of  $\mathbf{b}$ ; but if you do, you must explain why this assumption

is reasonable.

*Hint:* Use the duality theorem to see how  $z(\mathbf{b})$  depends on  $\mathbf{b}$ .

- b. Suppose that the first component of the optimal dual vector is  $\lambda_1 = 3$ . Now suppose that we increase  $b_1$  by a very small amount  $\Delta b_1$ . Determine the amount by which the optimal objective function value will change.

**17.24** Consider the *quadratic programming* problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ . Call this problem the *primal problem*.

Consider the associated *dual* quadratic programming problem

$$\begin{aligned} & \text{maximize} && -\frac{1}{2} \mathbf{y}^\top (\mathbf{A}\mathbf{A}^\top) \mathbf{y} - \mathbf{b}^\top \mathbf{y} \\ & \text{subject to} && \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

Let  $f_1$  and  $f_2$  be the objective functions of the primal and dual, respectively.

- a. State and prove a weak duality lemma in this setting.
- b. Show that if  $\mathbf{x}_0$  and  $\mathbf{y}_0$  are feasible points in the primal and dual, and  $f_1(\mathbf{x}_0) = f_2(\mathbf{y}_0)$ , then  $\mathbf{x}_0$  and  $\mathbf{y}_0$  are optimal solutions to the primal and dual, respectively.

*Hint:* The techniques used in the linear programming duality results are applicable in this exercise.



# CHAPTER 18

---

## NONSIMPLEX METHODS

---

### 18.1 Introduction

In previous chapters we studied the simplex method and its variant, the revised simplex method, for solving linear programming problems. The method remains widely used in practice for solving LP problems. However, the amount of time required to compute a solution using the simplex method grows rapidly as the number of components  $n$  of the variable  $\mathbf{x} \in \mathbb{R}^n$  increases. Specifically, it turns out that the relationship between the required amount of time for the algorithm to find a solution and the size  $n$  of  $\mathbf{x}$  is exponential in the worst case. An example of an LP problem for which this relationship is evident was devised by Klee and Minty in 1972 [76]. Below, we give a version of the

Klee-Minty example, taken from [9]. Let  $n$  be given. Let

$$\begin{aligned}\mathbf{c} &= [10^{n-1}, 10^{n-2}, \dots, 10^1, 1]^\top, \\ \mathbf{b} &= [1, 10^2, 10^4, \dots, 10^{2(n-1)}]^\top, \\ \mathbf{A} &= \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 2 \times 10^1 & 1 & 0 & \cdots & 0 \\ 2 \times 10^2 & 2 \times 10^1 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 2 \times 10^{n-1} & 2 \times 10^{n-2} & \cdots & 2 \times 10^1 & 1 \end{bmatrix}.\end{aligned}$$

Consider the following LP problem:

$$\begin{aligned}&\text{maximize} \quad \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ &\quad \mathbf{x} \geq \mathbf{0}.\end{aligned}$$

The simplex algorithm applied to the LP problem above requires  $2^n - 1$  steps to find the solution. Clearly, in this example the relationship between the required amount of time for the simplex algorithm to find a solution and the size  $n$  of the variable  $\mathbf{x}$  is exponential. This relationship is also called the *complexity* of the algorithm. The simplex algorithm is therefore said to have *exponential complexity*. The complexity of the simplex algorithm is also often written as  $O(2^n - 1)$ .

Naturally, we would expect that any algorithm that solves LP problems would have the property that the time required to arrive at a solution increases with the size  $n$  of the variable  $\mathbf{x}$ . However, the issue at hand is the rate at which this increase occurs. As we have seen above, the simplex algorithm has the property that this rate of increase is exponential. For a number of years, computer scientists have distinguished between *exponential complexity* and *polynomial complexity*. If an algorithm for solving LP problems has polynomial complexity, then the time required to obtain the solution is bounded by a polynomial in  $n$ . Obviously, polynomial complexity is more desirable than exponential complexity. Therefore, the existence of an algorithm for solving LP problems with polynomial complexity is an important issue. This issue was partially resolved in 1979 by Khachiyan (also transliterated as Hacijan) [74], who proposed an algorithm that has a complexity  $O(n^4 L)$ , where, roughly speaking,  $L$  represents the number of bits used in the computations. The reason that we consider Khachiyan's algorithm (also called the *ellipsoid algorithm*) as only a partial resolution of this issue is that the complexity depends on  $L$ , which implies that the time required to solve a given LP problem increases with the required accuracy of the computations. The existence of a method for solving LP problems with a polynomial complexity bound based only on the size of the variable  $n$  (and possibly the number of constraints) remains a difficult open problem [55]. In any case, computational

experience with Khachiyan's algorithm has shown that it is not a practical alternative to the simplex method [14]. The theoretical complexity advantage of Khachiyan's method relative to the simplex method remains to be demonstrated in practice.

Another nonsimplex algorithm for solving LP problems was proposed in 1984 by Karmarkar [71]. Karmarkar's algorithm has a complexity of  $O(n^{3.5}L)$ , which is lower than that of Khachiyan's algorithm. The algorithm is superior to the simplex algorithm from a complexity viewpoint, but has its drawbacks. Improved methods along similar lines, called *interior-point methods*, have received considerable interest since Karmarkar's original paper. Well-implemented versions of these methods are very efficient, especially when the problem involves a large number of variables [55].

This chapter is devoted to a discussion of nonsimplex methods for solving LP problems. In the next section we discuss some ideas underlying Khachiyan's algorithm. We then present Karmarkar's algorithm in the section to follow.

## 18.2 Khachiyan's Method

Our description of the Khachiyan's algorithm is based on [8] and [9]. The method relies on the concept of duality (see Chapter 17). Our exposition of Khachiyan's algorithm is geared toward a basic understanding of the method. For a detailed rigorous treatment of the method, we refer the reader to [101].

Consider the (primal) linear programming problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

We write the corresponding dual problem,

$$\begin{aligned} & \text{maximize} && \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to} && \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & && \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Recall that the two LP problems above constitute the symmetric form of duality. From Theorem 17.1, if  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are feasible solutions to the primal and dual problems, respectively, and  $\mathbf{c}^\top \mathbf{x} = \boldsymbol{\lambda}^\top \mathbf{b}$ , then  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  are optimal solutions to their respective problems. Using this result, we see that to solve the primal problem it is enough to find a vector  $[\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top$  that satisfies the

following set of relations:

$$\begin{aligned} \mathbf{c}^\top \mathbf{x} &= \mathbf{b}^\top \boldsymbol{\lambda}, \\ \mathbf{A}\mathbf{x} &\geq \mathbf{b}, \\ \mathbf{A}^\top \boldsymbol{\lambda} &\leq \mathbf{c}, \\ \mathbf{x} &\geq \mathbf{0}, \\ \boldsymbol{\lambda} &\geq \mathbf{0}. \end{aligned}$$

Note that the equality  $\mathbf{c}^\top \mathbf{x} = \mathbf{b}^\top \boldsymbol{\lambda}$  is equivalent to the two inequalities

$$\begin{aligned} \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\lambda} &\leq 0, \\ -\mathbf{c}^\top \mathbf{x} + \mathbf{b}^\top \boldsymbol{\lambda} &\leq 0. \end{aligned}$$

Taking this into account, we can represent the previous set of relations as

$$\begin{bmatrix} \mathbf{c}^\top & -\mathbf{b}^\top \\ -\mathbf{c}^\top & \mathbf{b}^\top \\ -\mathbf{A} & \mathbf{0} \\ -\mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{A}^\top \\ \mathbf{0} & -\mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} \leq \begin{bmatrix} 0 \\ 0 \\ -\mathbf{b} \\ \mathbf{0} \\ \mathbf{c} \\ \mathbf{0} \end{bmatrix}.$$

Therefore, we have reduced the problem of finding an optimal solution to the primal-dual pair into one of finding a vector  $[\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top$  that satisfies the system of inequalities above. In other words, if we can find a vector that satisfies the system of inequalities, then this vector gives an optimal solution to the primal-dual pair. On the other hand, if there does not exist a vector satisfying the system of inequalities, then the primal-dual pair has no optimal feasible solution. In the subsequent discussion, we simply represent the system of inequalities as

$$\mathbf{P}\mathbf{z} \leq \mathbf{q},$$

where

$$\mathbf{P} = \begin{bmatrix} \mathbf{c}^\top & -\mathbf{b}^\top \\ -\mathbf{c}^\top & \mathbf{b}^\top \\ -\mathbf{A} & \mathbf{0} \\ -\mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{A}^\top \\ \mathbf{0} & -\mathbf{I}_m \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix}, \quad \mathbf{q} = \begin{bmatrix} 0 \\ 0 \\ -\mathbf{b} \\ \mathbf{0} \\ \mathbf{c} \\ \mathbf{0} \end{bmatrix}.$$

In our discussion of Khachiyan's algorithm, we will not be using these forms of  $\mathbf{P}$ ,  $\mathbf{q}$ , and  $\mathbf{z}$  specifically; we simply treat  $\mathbf{P}\mathbf{z} \leq \mathbf{q}$  as a generic matrix inequality, with  $\mathbf{P}$ ,  $\mathbf{q}$ , and  $\mathbf{z}$  as generic entities. Let  $r$  and  $s$  be the sizes of  $\mathbf{q}$  and  $\mathbf{z}$ , respectively; that is,  $\mathbf{P} \in \mathbb{R}^{r \times s}$ ,  $\mathbf{z} \in \mathbb{R}^s$ , and  $\mathbf{q} \in \mathbb{R}^r$ .

Khachiyan's method solves the LP problem by first determining if there exists a vector  $\mathbf{z}$  that satisfies the inequality  $\mathbf{P}\mathbf{z} \leq \mathbf{q}$ ; that is, the algorithm decides if the system of linear inequalities above is *consistent*. If the system is consistent, then the algorithm finds a vector  $\mathbf{z}$  satisfying the system. In the following we refer to any vector satisfying the system as a *solution* to the system. We assume that the entries in  $\mathbf{P}$  and  $\mathbf{q}$  are all rational numbers. This is not a restriction in practice, since any representation of our LP problem on a digital computer will involve only rational numbers. In fact, we assume further that the entries in  $\mathbf{P}$  and  $\mathbf{q}$  are all integers. We can do this without loss of generality since we can always multiply both sides of the inequality  $\mathbf{P}\mathbf{z} \leq \mathbf{q}$  by a sufficiently large number to get only integer entries on both sides.

Before discussing Khachiyan's algorithm, we introduce the idea of an *ellipsoid*. To this end, let  $\mathbf{z} \in \mathbb{R}^s$  be a given vector and let  $\mathbf{Q}$  be an  $s \times s$  nonsingular matrix. Then, the *ellipsoid* associated with  $\mathbf{Q}$  centered at  $\mathbf{z}$  is defined as the set

$$E_{\mathbf{Q}}(\mathbf{z}) = \{\mathbf{z} + \mathbf{Q}\mathbf{y} : \mathbf{y} \in \mathbb{R}^s, \|\mathbf{y}\| \leq 1\}.$$

The main idea underlying Khachiyan's algorithm is as follows. Khachiyan's algorithm is an iterative procedure, where at each iteration we update a vector  $\mathbf{z}^{(k)}$  and a matrix  $\mathbf{Q}_k$ . Associated with  $\mathbf{z}^{(k)}$  and  $\mathbf{Q}_k$  is an ellipsoid  $E_{\mathbf{Q}_k}(\mathbf{z}^{(k)})$ . At each step of the algorithm, the associated ellipsoid contains a solution to the given system of linear inequalities. The algorithm updates  $\mathbf{z}^{(k)}$  and  $\mathbf{Q}_k$  in such a way that the ellipsoid at the next step is "smaller" than that of the current step, but at the same time is guaranteed to contain a solution to the given system of inequalities, if one exists. If we find that the current point  $\mathbf{z}^{(k)}$  satisfies  $\mathbf{P}\mathbf{z}^{(k)} \leq \mathbf{q}$ , then we terminate the algorithm and conclude that  $\mathbf{z}^{(k)}$  is a solution. Otherwise, we continue to iterate. The algorithm has a fixed prespecified maximum number of iterations  $N$  to be performed, where  $N$  is a number that depends on  $L$  and  $s$ . Note that we are not free to choose  $N$ —it is computed using a formula that uses the values of  $L$  and  $s$ . The constant  $L$  is itself a quantity that we have to compute beforehand using a formula that involves  $\mathbf{P}$  and  $\mathbf{q}$ . When we have completed  $N$  iterations without finding a solution in an earlier step, we terminate the algorithm. The associated ellipsoid will then have shrunk to the extent that it is smaller than the precision of computation. At this stage, we will either discover a solution inside the ellipsoid, if indeed a solution exists, or we will find no solution inside the ellipsoid, in which case we conclude that no solution exists.

As we can see from the description above, Khachiyan's approach is a radical departure from the classical simplex method for solving LP problems. The method has attracted a lot of attention, and many studies have been devoted to it. However, as we pointed out earlier, the algorithm is of little practical value for solving real-world LP problems. Therefore, we do not delve any further into the details of Khachiyan's algorithm. We refer the interested reader to [101].

Despite its practical drawbacks, Khachiyan's method has inspired other researchers to pursue the development of computationally efficient algorithms for solving LP problems with polynomial complexity. One such algorithm is attributed to Karmarkar, which we discuss in Section 18.4.

### 18.3 Affine Scaling Method

#### Basic Algorithm

In this section we describe a simple algorithm, called the *affine scaling method*, for solving linear programming problems. This description is to prepare the reader for our discussion of Karmarkar's method in the next section. The affine scaling method is a *an interior-point method*. Such methods differ fundamentally from the classical simplex method in one main respect: An interior-point method starts inside the feasible set and moves within it toward an optimal vertex. In contrast, the simplex method jumps from vertex to vertex of the feasible set seeking an optimal vertex.

To begin our description of the affine scaling method, consider the LP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

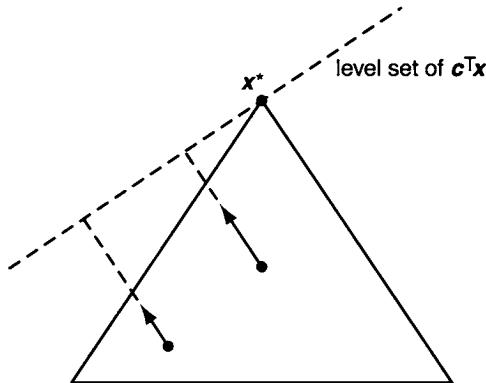
Note that the feasibility constraints have two parts:  $\mathbf{A}\mathbf{x} = \mathbf{b}$  and  $\mathbf{x} \geq \mathbf{0}$ . Suppose that we have a feasible point  $\mathbf{x}^{(0)}$  that is *strictly interior* (by strictly interior we mean that all of the components of  $\mathbf{x}^{(0)}$  are strictly positive). We wish to find a new point  $\mathbf{x}^{(1)}$  by searching in a direction  $\mathbf{d}^{(0)}$  that decreases the objective function. In other words, we set

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)},$$

where  $\alpha_0$  is a step size. In the gradient method (Chapter 8) we used the negative gradient of the objective function for the search direction. For the LP problem, the negative gradient of the objective function is  $-\mathbf{c}$ . However, if we set  $\mathbf{d}^{(0)} = -\mathbf{c}$ , the point  $\mathbf{x}^{(1)}$  may not lie inside the feasible set. For  $\mathbf{x}^{(1)}$  to lie inside the feasible set, it is necessary that  $\mathbf{d}^{(0)}$  be a vector in the nullspace of  $\mathbf{A}$ . Indeed, because  $\mathbf{x}^{(0)}$  is feasible, we have  $\mathbf{A}\mathbf{x}^{(0)} = \mathbf{b}$ . We also require that  $\mathbf{A}\mathbf{x}^{(1)} = \mathbf{b}$ . Combining these two equations yields

$$\mathbf{A}(\mathbf{x}^{(1)} - \mathbf{x}^{(0)}) = \alpha_0 \mathbf{A}\mathbf{d}^{(0)} = \mathbf{0}.$$

To choose a direction  $\mathbf{d}^{(0)}$  that lies in the nullspace of  $\mathbf{A}$  but is still “close” to  $-\mathbf{c}$ , we orthogonally project  $-\mathbf{c}$  onto the nullspace of  $\mathbf{A}$  and take the resulting



**Figure 18.1** Results of projected gradient step from center and noncenter points.

projection as  $\mathbf{d}^{(0)}$ . The orthogonal projection of any vector onto the nullspace of  $\mathbf{A}$  involves multiplication by the following matrix  $\mathbf{P}$ , called the *orthogonal projector* (see Section 3.3 and Example 12.5):

$$\mathbf{P} = \mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}.$$

We set  $\mathbf{d}^{(0)}$  to be in the direction of the orthogonal projection of  $-\mathbf{c}$  onto the nullspace of  $\mathbf{A}$ :

$$\mathbf{d}^{(0)} = -\mathbf{P}\mathbf{c}.$$

It is easy to check that  $\mathbf{A}\mathbf{P}\mathbf{c} = \mathbf{0}$  and hence  $\mathbf{A}\mathbf{x}^{(1)} = \mathbf{b}$ . In summary, given a feasible point  $\mathbf{x}^{(0)}$ , we find a new feasible point  $\mathbf{x}^{(1)}$  using

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \mathbf{P}\mathbf{c},$$

where the choice of the step size  $\alpha_0$  is discussed later in the section. The choice of  $\mathbf{x}^{(1)}$  above can be viewed as one iteration of a *projected gradient algorithm*, discussed in Section 23.3.

We now make the observation that the point  $\mathbf{x}^{(0)}$  should be chosen close to the center of the feasible set. Figure 18.1 illustrates this observation. Comparing the center and noncenter starting points in the figure, we can see that if we start at the center of the feasible set, we can take a larger step in the search direction. This larger step from the center point should yield a lower-cost value for the new point compared with the step originating from the noncenter point.

Suppose that we are given a point  $\mathbf{x}^{(0)}$  that is feasible but is not a center point. We can transform the point to the center by applying what is called an *affine scaling*. For simplicity, suppose that  $\mathbf{A} = [1, 1, \dots, 1]/n$  and  $\mathbf{b} = [1]$ . It is easy to see that the center of this feasible set is the point  $\mathbf{e} = [1, \dots, 1]^\top$ . To transform  $\mathbf{x}^{(0)}$  to  $\mathbf{e}$ , we use the affine-scaling transformation

$$\mathbf{e} = \mathbf{D}_0^{-1} \mathbf{x}^{(0)},$$

where  $\mathbf{D}_0$  is a diagonal matrix whose diagonal entries are the components of the vector  $\mathbf{x}^{(0)}$ :

$$\mathbf{D}_0 = \text{diag}[x_1^{(0)}, \dots, x_n^{(0)}] = \begin{bmatrix} x_1^{(0)} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & x_n^{(0)} \end{bmatrix}.$$

Note that  $\mathbf{D}_0$  is invertible because we assumed that  $\mathbf{x}^{(0)}$  is *strictly interior*. For general  $\mathbf{A}$  and  $\mathbf{b}$  we will still use the same affine-scaling transformation as above. In general, we may not be at precisely the center of the feasible set, but we hope that the transformed point will be “close” to the center. At least the point  $\mathbf{e}$  is equidistant from the boundaries of the positive orthant  $\{\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ .

Once the starting point is at (or close to) the center of the feasible set after performing the affine-scaling transformation, we can proceed as described before. Because we have transformed the original vector  $\mathbf{x}^{(0)}$  via premultiplication by  $\mathbf{D}_0^{-1}$ , effectively changing the coordinate system, we also need to represent the original LP problem in the new coordinates. Specifically, the LP problem in the transformed coordinates takes the form

$$\begin{aligned} & \text{minimize} && \bar{\mathbf{c}}_0^\top \bar{\mathbf{x}} \\ & \text{subject to} && \bar{\mathbf{A}}_0 \bar{\mathbf{x}} = \mathbf{b} \\ & && \bar{\mathbf{x}} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \bar{\mathbf{c}}_0 &= \mathbf{D}_0 \mathbf{c}, \\ \bar{\mathbf{A}}_0 &= \mathbf{A} \mathbf{D}_0. \end{aligned}$$

In the new  $(\bar{\mathbf{x}})$  coordinate system we construct the orthogonal projector

$$\bar{\mathbf{P}}_0 = \mathbf{I}_n - \bar{\mathbf{A}}_0^\top (\bar{\mathbf{A}}_0 \bar{\mathbf{A}}_0^\top)^{-1} \bar{\mathbf{A}}_0$$

and set  $\bar{\mathbf{d}}^{(0)}$  to be in the direction of the orthogonal projection of  $-\bar{\mathbf{c}}_0$  onto the nullspace of  $\bar{\mathbf{A}}_0$ :

$$\bar{\mathbf{d}}^{(0)} = -\bar{\mathbf{P}}_0 \bar{\mathbf{c}}_0.$$

Then, compute  $\bar{\mathbf{x}}^{(1)}$  using

$$\bar{\mathbf{x}}^{(1)} = \bar{\mathbf{x}}^{(0)} - \alpha_0 \bar{\mathbf{P}}_0 \bar{\mathbf{c}}_0,$$

where  $\bar{\mathbf{x}}^{(0)} = \mathbf{D}_0^{-1} \mathbf{x}^{(0)}$ . To obtain a point in the original coordinates, we perform the transformation

$$\mathbf{x}^{(1)} = \mathbf{D}_0 \bar{\mathbf{x}}^{(1)}.$$

The procedure above takes a point  $\mathbf{x}^{(0)}$  and generates a new point  $\mathbf{x}^{(1)}$ . This procedure can be represented as

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)},$$

where

$$\mathbf{d}^{(0)} = -\mathbf{D}_0 \bar{\mathbf{P}} \mathbf{D}_0 \mathbf{c}.$$

We repeat the procedure iteratively to generate a sequence of points  $\{\mathbf{x}^{(k)}\}$ , where

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$$

with

$$\begin{aligned}\mathbf{D}_k &= \text{diag}[x_1^{(k)}, \dots, x_n^{(k)}], \\ \bar{\mathbf{A}}_k &= \mathbf{A} \mathbf{D}_k, \\ \bar{\mathbf{P}}_k &= \mathbf{I}_n - \bar{\mathbf{A}}_k^\top (\bar{\mathbf{A}}_k \bar{\mathbf{A}}_k^\top)^{-1} \bar{\mathbf{A}}_k, \\ \mathbf{d}^{(k)} &= -\mathbf{D}_k \bar{\mathbf{P}}_k \mathbf{D}_k \mathbf{c}.\end{aligned}$$

At each stage of the algorithm, we have to ensure that the point  $\mathbf{x}^{(k)}$  is strictly interior. Note that the condition  $\mathbf{A}\mathbf{x}^{(k)} = \mathbf{b}$  is satisfied automatically at each stage because of the way we select  $\mathbf{d}^{(k)}$ . However, we also need to guarantee that  $x_i^{(k)} > 0$  for  $i = 1, \dots, n$ . This can be done through appropriate choice of the step size  $\alpha_k$ , discussed next.

The main criterion for choosing  $\alpha_k$  is to make it as large as possible, but not so large that some components of  $\mathbf{x}^{(k+1)}$  become nonpositive. That is, we select  $\alpha_k$  so that  $x_i^{(k+1)} = x_i^{(k)} + \alpha_k d_i^{(k)} > 0$  for  $i = 1, \dots, n$ . To proceed, first define

$$r_k = \min_{\{i: d_i^{(k)} < 0\}} -\frac{x_i^{(k)}}{d_i^{(k)}}.$$

The number  $r_k$  represents the largest value of the step size  $\alpha_k$  such that all the components of  $\mathbf{x}^{(k+1)}$  are nonnegative. To ensure that  $\mathbf{x}^{(k+1)}$  is strictly interior, we use a step size of the form  $\alpha_k = \alpha r_k$ , where  $\alpha \in (0, 1)$ . Typical values of  $\alpha$  for this method are  $\alpha = 0.9$  or  $0.99$  (see [96, p. 572]).

Unlike the simplex method, the affine scaling method will not reach the optimal solution in a finite number of steps. Therefore, we need a stopping criterion. For this, we can use any of the stopping criteria discussed in Section 8.2. For example, we can stop if

$$\frac{|\mathbf{c}\mathbf{x}^{(k+1)} - \mathbf{c}\mathbf{x}^{(k)}|}{\max\{1, |\mathbf{c}\mathbf{x}^{(k)}|\}} < \varepsilon,$$

where  $\varepsilon > 0$  is a prespecified threshold (see also [96, p. 572] for a similar stopping criterion, as well as an alternative criterion involving duality).

## Two-Phase Method

To implement the affine scaling method described above, we need an initial feasible starting point that is strictly interior. We now describe a method to find such a starting point. After the starting point is found, we can then proceed to search for an optimal solution to the problem. This approach involves two phases: In phase I we find an initial strictly interior feasible point, and in phase II we use the result of phase I to initialize the affine scaling algorithm to find an optimal solution. This procedure is analogous to the two-phase simplex algorithm described in Section 16.6.

We now describe phase I of the two-phase affine scaling method. Let  $\mathbf{u}$  be an arbitrary vector with positive components, and let

$$\mathbf{v} = \mathbf{b} - \mathbf{A}\mathbf{u}.$$

If  $\mathbf{v} = \mathbf{0}$ , then  $\mathbf{u}$  is a strictly interior feasible point. We can then set  $\mathbf{x}^{(0)} = \mathbf{u}$  and proceed to phase II, where we apply the affine scaling method as described before. On the other hand, if  $\mathbf{v} \neq \mathbf{0}$ , we construct the following associated *artificial problem*:

$$\begin{aligned} & \text{minimize} && y \\ & \text{subject to} && [\mathbf{A}, \mathbf{v}] \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} = \mathbf{b} \\ & && \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \geq \mathbf{0}. \end{aligned}$$

The artificial problem above has an obvious strictly interior feasible point:

$$\begin{bmatrix} \mathbf{u} \\ 1 \end{bmatrix}.$$

Using this point as the initial point, we can apply the affine scaling algorithm to the artificial problem. Because the objective function in the artificial problem is bounded below by 0, the affine scaling method will terminate with some optimal solution.

**Proposition 18.1** *The original LP problem has a feasible solution if and only if the associated artificial problem has an optimal feasible solution with objective function value zero.*  $\square$

*Proof.*  $\Rightarrow$ : If the original problem has a feasible solution  $\mathbf{x}$ , then the vector  $[\mathbf{x}^\top, 0]^\top$  is a feasible solution to the artificial problem. Clearly, this solution has an objective function value of zero. This solution is therefore optimal for the artificial problem, since there can be no feasible solution with negative objective function value.

$\Leftarrow$ : Suppose that the artificial problem has an optimal feasible solution with objective function value zero. Then, this solution must have the form  $[\mathbf{x}^\top, 0]^\top$ , where  $\mathbf{x} \geq \mathbf{0}$ . Hence, we have  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , and  $\mathbf{x}$  is a feasible solution to the original problem. ■

Suppose that the original LP problem has a feasible solution. By Proposition 18.1, if we apply the affine scaling method to the artificial problem (with initial point  $[\mathbf{u}^\top, 1]^\top$ ), the algorithm will terminate with objective function value zero. The optimal solution will be of the form  $[\mathbf{x}^\top, 0]^\top$ . We argue that  $\mathbf{x}$  will in general be a strictly interior feasible point. It is easy to see that  $\mathbf{x} \geq \mathbf{0}$ . To convince ourselves that each component of  $\mathbf{x}$  will be positive in general, note that the subset of optimal feasible solutions of the artificial problem in which one or more among the first  $n$  components are zero is a very small or thin subset of the set of all optimal feasible solutions. By *small* or *thin* we mean in the sense that a two-dimensional plane in  $\mathbb{R}^3$  is small or thin. In particular, the volume of the two-dimensional plane in  $\mathbb{R}^3$  is zero. Thus, it is very unlikely that the affine scaling algorithm will terminate with an optimal feasible solution in which one or more among the first  $n$  components are zero.

Having completed phase I as described above, we then use the first  $n$  components of the terminal optimal feasible solution for the artificial problem as our initial point for the affine scaling method applied to the original LP problem. This second application of the affine scaling algorithm constitutes phase II.

In theory, phase I generates a feasible point to initiate phase II. However, because of the finite precision of typical computer implementations, the solution obtained from phase I may not, in fact, be feasible. Moreover, even if the initial point in phase II is feasible, in practice the iterates may lose feasibility, owing to finite precision computations. Special procedures for dealing with such problems are available. For a discussion of numerical implementation of affine scaling algorithms, see [42, Section 7.1.2].

## 18.4 Karmarkar's Method

### Basic Ideas

Like the affine scaling method, Karmarkar's method for solving LP problems differs fundamentally from the classical simplex method in various respects. First, Karmarkar's method is an interior-point method. Another difference between Karmarkar's method and the simplex method is that the latter stops when it finds an optimal solution. On the other hand, Karmarkar's method stops when it finds a solution that has an objective function value that is less than or equal to a prespecified fraction of the original guess. A third difference between the two methods is that the simplex method starts with LP problems in standard form, whereas Karmarkar's method starts with LP problems in a special canonical form, which we call *Karmarkar's canonical*

*form.* We discuss this canonical form in the next subsection. While more recent interior-point methods are recognized to be superior to Karmarkar's original algorithm in efficiency and robustness, a study of Karmarkar's method provides an informative introduction to the study of more advanced interior-point methods.

### Karmarkar's Canonical Form

To apply Karmarkar's algorithm to a given LP problem, we must first transform the given problem into a particular form, which we refer to as Karmarkar's canonical form. Karmarkar's canonical form is written as

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{0} \\ & && \sum_{i=1}^n x_i = 1 \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{x} = [x_1, \dots, x_n]^\top$ . As in our discussion of Khachiyan's method, we assume without loss of generality that the entries of  $\mathbf{A}$  and  $\mathbf{c}$  are integers.

We now introduce some notation that allows convenient manipulation of the canonical form. First, let  $\mathbf{e} = [1, \dots, 1]^\top$  be the vector in  $\mathbb{R}^n$  with each component equal to 1. Let  $\Omega$  denote the nullspace of  $\mathbf{A}$ , that is, the subspace

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}.$$

Define the *simplex* in  $\mathbb{R}^n$  by

$$\Delta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}.$$

We denote the *center* of the simplex  $\Delta$  by

$$\mathbf{a}_0 = \frac{\mathbf{e}}{n} = \left[ \frac{1}{n}, \dots, \frac{1}{n} \right]^\top.$$

Clearly,  $\mathbf{a}_0 \in \Delta$ . With the notation above, Karmarkar's canonical form can be rewritten as

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{x} \in \Omega \cap \Delta. \end{aligned}$$

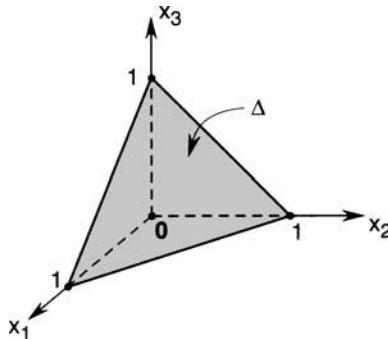
Note that the constraint set (or feasible set)  $\Omega \cap \Delta$  can be represented as

$$\begin{aligned} \Omega \cap \Delta &= \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}, \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^n : \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \mathbf{x} \geq \mathbf{0} \right\}. \end{aligned}$$

**Example 18.1** Consider the following LP problem, taken from [125]:

$$\begin{aligned} & \text{minimize} && 5x_1 + 4x_2 + 8x_3 \\ & \text{subject to} && x_1 + x_2 + x_3 = 1 \\ & && x_1, x_2, x_3 \geq 0. \end{aligned}$$

Clearly, this problem is already in Karmarkar's canonical form, with  $\mathbf{c}^\top = [5, 4, 8]$ , and  $\mathbf{A} = \mathbf{O}$ . The feasible set for this example is illustrated in Figure 18.2. ■



**Figure 18.2** Feasible set for Example 18.1.

**Example 18.2** Consider the following LP problem, taken from [110]:

$$\begin{aligned} & \text{minimize} && 3x_1 + 3x_2 - x_3 \\ & \text{subject to} && 2x_1 - 3x_2 + x_3 = 0 \\ & && x_1 + x_2 + x_3 = 1 \\ & && x_1, x_2, x_3 \geq 0. \end{aligned}$$

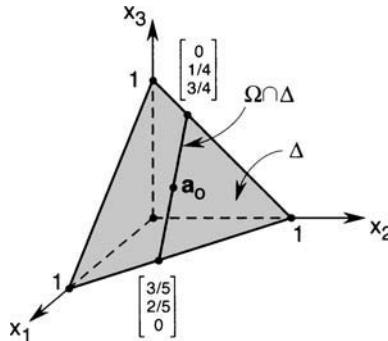
This problem is in Karmarkar's canonical form, with  $\mathbf{c}^\top = [3, 3, -1]$  and  $\mathbf{A} = [2, -3, 1]$ . The feasible set for this example is illustrated in Figure 18.3 (adapted from [110]). ■

We show later that any LP problem can be converted into an equivalent problem in Karmarkar's canonical form.

### Karmarkar's Restricted Problem

Karmarkar's algorithm solves LP problems in Karmarkar's canonical form, with the following assumptions:

- A. The center  $\mathbf{a}_0$  of the simplex  $\Delta$  is a feasible point:  $\mathbf{a}_0 \in \Omega$ .



**Figure 18.3** The feasible set for Example 18.2.

B. The minimum value of the objective function over the feasible set is zero.

C. The  $(m + 1) \times n$  matrix

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix}$$

has rank  $m + 1$ .

D. We are given a termination parameter  $q > 0$ , such that if we obtain a feasible point  $\mathbf{x}$  satisfying

$$\frac{\mathbf{c}^\top \mathbf{x}}{\mathbf{c}^\top \mathbf{a}_0} \leq 2^{-q},$$

then we consider the problem solved.

Any LP problem that is in Karmarkar's canonical form and that also satisfies the four assumptions above is called a *Karmarkar's restricted problem*. In the following we discuss the assumptions and their interpretations.

We begin by looking at assumption A. We point out that this assumption is not restrictive, since any LP problem that has an optimal feasible solution can be converted into a problem in Karmarkar's canonical form that satisfies assumption A. We discuss this in the next subsection.

We next turn our attention to assumption B. Any LP problem in Karmarkar's canonical form can be converted into one that satisfies assumption B, provided that we know beforehand the minimum value of its objective function over the feasible set. Specifically, suppose that we are given an LP problem where the minimum value of the objective function is  $M$ . As in [110], consider the function  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x} - M$ . Then, using the property that  $\mathbf{e}^\top \mathbf{x} = 1$  on the feasible set, we have that for any feasible  $\mathbf{x}$ ,

$$f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x} - M = \mathbf{c}^\top \mathbf{x} - M \mathbf{e}^\top \mathbf{x} = (\mathbf{c}^\top - M \mathbf{e}^\top) \mathbf{x} = \tilde{\mathbf{c}}^\top \mathbf{x},$$

where  $\tilde{\mathbf{c}}^\top = \mathbf{c}^\top - M\mathbf{e}^\top$ . Notice that the objective function above has a minimum value of zero and is a linear function of  $\mathbf{x}$ . We can replace the original objective function with the new objective function above, without altering the solution.

**Example 18.3** Recall the LP problem in Example 18.1:

$$\begin{aligned} & \text{minimize} && 5x_1 + 4x_2 + 8x_3 \\ & \text{subject to} && x_1 + x_2 + x_3 = 1 \\ & && x_1, x_2, x_3 \geq 0. \end{aligned}$$

The problem satisfies assumption A (and assumption C) but not assumption B, since the minimum value of the objective function over the feasible set is 4. To convert the above into a problem that satisfies assumption B, we replace  $\mathbf{c}^\top = [5, 4, 8]$  by  $\tilde{\mathbf{c}}^\top = [1, 0, 4]$ . ■

**Example 18.4** The reader can easily verify that the LP problem in Example 18.2 satisfies assumptions A, B, and C. ■

Assumption C is a technical assumption that is required in the implementation of the algorithm. Its significance will be clear when we discuss the update equation in Karmarkar's algorithm.

Assumption D is the basis for the stopping criterion of Karmarkar's algorithm. In particular, we stop when we have found a feasible point satisfying  $\mathbf{c}^\top \mathbf{x}/\mathbf{c}^\top \mathbf{a}_0 \leq 2^{-q}$ . Such a stopping criterion is inherent in any algorithm that uses finite-precision arithmetic. Observe that the stopping criterion above depends on the value of  $\mathbf{c}^\top \mathbf{a}_0$ . It will turn out that Karmarkar's algorithm uses  $\mathbf{a}_0$  as the starting point. Therefore, we can see that the accuracy of the final solution in the algorithm is influenced by the starting point.

### From General Form to Karmarkar's Canonical Form

We now show how any LP problem can be converted into an equivalent problem in Karmarkar's canonical form. By *equivalent* we mean that the solution to one can be used to determine the solution to the other, and vice versa. To this end, recall that any LP problem can be transformed into an equivalent problem in standard form. Therefore, it suffices to show that any LP problem in standard form can be transformed into an equivalent problem in Karmarkar's canonical form. In fact, the transformation given below (taken from [71]) will also guarantee that assumption A of the preceding subsection is satisfied.

To proceed, consider a given LP problem in standard form:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n \\ & \text{subject to} && \mathbf{Ax} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

We first present a simple way to convert this problem into Karmarkar's canonical form, ignoring the requirement to satisfy assumption A. For this, define a new variable  $\mathbf{z} \in \mathbb{R}^{n+1}$  by

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}.$$

Also define  $\mathbf{c}' = [\mathbf{c}^\top, 0]^\top$  and  $\mathbf{A}' = [\mathbf{A}, -\mathbf{b}]$ . Using this notation, we can now rewrite the LP problem above as

$$\begin{aligned} & \text{minimize} && \mathbf{c}'^\top \mathbf{z}, \quad \mathbf{z} \in \mathbb{R}^{n+1} \\ & \text{subject to} && \mathbf{A}' \mathbf{z} = \mathbf{0} \\ & && \mathbf{z} \geq \mathbf{0}. \end{aligned}$$

We need one more step to transform the problem into one that includes the constraint that the decision variables sum to 1. For this, let  $\mathbf{y} = [y_1, \dots, y_n, y_{n+1}]^\top \in \mathbb{R}^{n+1}$ , where

$$\begin{aligned} y_i &= \frac{x_i}{x_1 + \dots + x_n + 1}, \quad i = 1, \dots, n \\ y_{n+1} &= \frac{1}{x_1 + \dots + x_n + 1}. \end{aligned}$$

This transformation from  $\mathbf{x}$  to  $\mathbf{y}$  is called a *projective transformation*. It can be shown that (see later)

$$\begin{aligned} \mathbf{c}^\top \mathbf{x} = 0 &\Leftrightarrow \mathbf{c}'^\top \mathbf{y} = 0, \\ \mathbf{A}\mathbf{x} = \mathbf{b} &\Leftrightarrow \mathbf{A}'\mathbf{y} = \mathbf{0}, \\ \mathbf{x} \geq \mathbf{0} &\Leftrightarrow \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

Therefore, we have transformed the given LP problem in standard form into the following problem, which is in Karmarkar's canonical form:

$$\begin{aligned} & \text{minimize} && \mathbf{c}'^\top \mathbf{y}, \quad \mathbf{y} \in \mathbb{R}^{n+1} \\ & \text{subject to} && \mathbf{A}' \mathbf{y} = \mathbf{0} \\ & && \mathbf{e}^\top \mathbf{y} = 1 \\ & && \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

The transformation technique above can be modified slightly to ensure that assumption A holds. We follow the treatment of [71]. We first assume that we are given a point  $\mathbf{a} = [a_1, \dots, a_n]$  that is a *strictly interior* feasible point; that is,  $\mathbf{A}\mathbf{a} = \mathbf{b}$  and  $\mathbf{a} > \mathbf{0}$ . We show later how this assumption can be enforced. Let  $P_+$  denote the *positive orthant* of  $\mathbb{R}^n$ , given by  $P_+ = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}\}$ . Let  $\Delta = \{\mathbf{x} \in \mathbb{R}^{n+1} : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$  be the simplex in  $\mathbb{R}^{n+1}$ . Define the map  $\mathbf{T} : P_+ \rightarrow \Delta$  by

$$\mathbf{T}(\mathbf{x}) = [T_1(\mathbf{x}), \dots, T_{n+1}(\mathbf{x})]^\top$$

with

$$\begin{aligned} T_i(\mathbf{x}) &= \frac{x_i/a_i}{x_1/a_1 + \cdots + x_n/a_n + 1}, \quad i = 1, \dots, n \\ T_{n+1}(\mathbf{x}) &= \frac{1}{x_1/a_1 + \cdots + x_n/a_n + 1}. \end{aligned}$$

We call the map  $\mathbf{T}$  a *projective transformation* of the positive orthant  $P_+$  into the simplex  $\Delta$  (for an introduction to projective transformations, see [68]). The transformation  $\mathbf{T}$  has several interesting properties (see Exercises 18.4, 18.5, and 18.6). In particular, we can find a vector  $\mathbf{c}' \in \mathbb{R}^{n+1}$  and a matrix  $\mathbf{A}' \in \mathbb{R}^{m \times (n+1)}$  such that for each  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\mathbf{c}'^\top \mathbf{x} = 0 \Leftrightarrow \mathbf{c}'^\top \mathbf{T}(\mathbf{x}) = 0$$

and

$$\mathbf{A}\mathbf{x} = \mathbf{b} \Leftrightarrow \mathbf{A}'\mathbf{T}(\mathbf{x}) = \mathbf{0}$$

(see Exercises 18.5 and 18.6 for the forms of  $\mathbf{A}'$  and  $\mathbf{c}'$ ). Note that for each  $\mathbf{x} \in \mathbb{R}^n$ , we have  $\mathbf{e}^\top \mathbf{T}(\mathbf{x}) = 1$ , which means that  $\mathbf{T}(\mathbf{x}) \in \Delta$ . Furthermore, note that for each  $\mathbf{x} \in \mathbb{R}^n$ ,

$$\mathbf{x} \geq 0 \Leftrightarrow \mathbf{T}(\mathbf{x}) \geq 0.$$

Taking this into account, consider the following LP problem (where  $\mathbf{y}$  is the decision variable):

$$\begin{aligned} &\text{minimize} && \mathbf{c}'^\top \mathbf{y} \\ &\text{subject to} && \mathbf{A}'\mathbf{y} = \mathbf{0} \\ &&& \mathbf{e}^\top \mathbf{y} = 1 \\ &&& \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

Note that this LP problem is in Karmarkar's canonical form. Furthermore, in light of the definitions of  $\mathbf{c}'$  and  $\mathbf{A}'$ , the above LP problem is equivalent to the original LP problem in standard form. Hence, we have converted the LP problem in standard form into an equivalent problem in Karmarkar's canonical form. In addition, because  $\mathbf{a}$  is a strictly interior feasible point, and  $\mathbf{a}_0 = \mathbf{T}(\mathbf{a})$  is the center of the simplex  $\Delta$  (see Exercise 18.4), the point  $\mathbf{a}_0$  is a feasible point of the transformed problem. Hence, assumption A of the preceding subsection is satisfied for the problem above.

We started this discussion with the assumption that we are given a point  $\mathbf{a}$  that is a strictly interior feasible point of the original LP problem in standard form. To see how this assumption can be made to hold, we now show that we can transform any given LP problem into an equivalent problem in standard form where such a point  $\mathbf{a}$  is explicitly given. To this end, consider a given

LP problem of the form

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & \quad \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Note that any LP problem can be converted into an equivalent problem of the above form. To see this, recall that any LP problem can be transformed into an equivalent problem in standard form. But any problem in standard form can be represented as above, since the constraint  $\mathbf{Ax} = \mathbf{b}$  can be written as  $\mathbf{Ax} \geq \mathbf{b}$ ,  $-\mathbf{Ax} \geq -\mathbf{b}$ . We next write the dual to the problem above:

$$\begin{aligned} & \text{maximize } \boldsymbol{\lambda}^\top \mathbf{b} \\ & \text{subject to } \boldsymbol{\lambda}^\top \mathbf{A} \leq \mathbf{c}^\top \\ & \quad \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

As we did in our discussion of Khachiyan's algorithm, we now combine the primal and dual problems to get

$$\begin{aligned} & \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\lambda} = 0, \\ & \mathbf{Ax} \geq \mathbf{b}, \\ & \mathbf{A}^\top \boldsymbol{\lambda} \leq \mathbf{c}, \\ & \mathbf{x} \geq \mathbf{0}, \\ & \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

As we pointed out in the earlier section on Khachiyan's algorithm, the original LP problem is solved if and only if we can find a pair  $(\mathbf{x}, \boldsymbol{\lambda})$  that satisfies the set of relations above. This follows from Theorem 17.1. We now introduce slack and surplus variables  $\mathbf{u}$  and  $\mathbf{v}$  to get the following equivalent set of relations:

$$\begin{aligned} & \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\lambda} = 0, \\ & \mathbf{Ax} - \mathbf{v} = \mathbf{b}, \\ & \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{u} = \mathbf{c}, \\ & \mathbf{x}, \boldsymbol{\lambda}, \mathbf{u}, \mathbf{v} \geq \mathbf{0}. \end{aligned}$$

Let  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\boldsymbol{\lambda}_0 \in \mathbb{R}^m$ ,  $\mathbf{u}_0 \in \mathbb{R}^n$ , and  $\mathbf{v}_0 \in \mathbb{R}^m$  be points that satisfy  $\mathbf{x}_0 > \mathbf{0}$ ,  $\boldsymbol{\lambda}_0 > \mathbf{0}$ ,  $\mathbf{u}_0 > \mathbf{0}$ , and  $\mathbf{v}_0 > \mathbf{0}$ . For example, we could choose  $\mathbf{x}_0 = [1, \dots, 1]^\top$ , and likewise with  $\boldsymbol{\lambda}_0$ ,  $\mathbf{u}_0$ , and  $\mathbf{v}_0$ . Consider the LP problem

$$\begin{aligned} & \text{minimize } z \\ & \text{subject to } \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\lambda} + (-\mathbf{c}^\top \mathbf{x}_0 + \mathbf{b}^\top \boldsymbol{\lambda}_0)z = 0 \\ & \quad \mathbf{Ax} - \mathbf{v} + (\mathbf{b} - \mathbf{Ax}_0 + \mathbf{v}_0)z = \mathbf{b} \\ & \quad \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{u} + (\mathbf{c} - \mathbf{A}^\top \boldsymbol{\lambda}_0)z = \mathbf{c} \\ & \quad \mathbf{x}, \boldsymbol{\lambda}, \mathbf{u}, \mathbf{v}, z \geq \mathbf{0}. \end{aligned}$$

We refer to the above as the *Karmarkar's artificial problem*, which can be represented in matrix notation as

$$\begin{aligned} & \text{minimize} && \tilde{\mathbf{c}}^\top \tilde{\mathbf{x}} \\ & \text{subject to} && \tilde{\mathbf{A}} \tilde{\mathbf{x}} = \tilde{\mathbf{b}} \\ & && \tilde{\mathbf{x}} \geq \mathbf{0}, \end{aligned}$$

where

$$\tilde{\mathbf{x}} = [\mathbf{x}^\top, \boldsymbol{\lambda}^\top, \mathbf{u}^\top, \mathbf{v}^\top, z]^\top,$$

$$\tilde{\mathbf{c}} = [\mathbf{0}_{2m+2n}^\top, 1]^\top,$$

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{c}^\top & -\mathbf{b}^\top & \mathbf{0}_n^\top & \mathbf{0}_m^\top & (-\mathbf{c}^\top \mathbf{x}_0 + \mathbf{b}^\top \boldsymbol{\lambda}_0) \\ \mathbf{A} & \mathbf{O}_{m \times m} & \mathbf{O}_{m \times n} & -\mathbf{I}_m & (\mathbf{b} - \mathbf{A}\mathbf{x}_0 + \mathbf{v}_0) \\ \mathbf{O}_{n \times n} & \mathbf{A}^\top & \mathbf{I}_n & \mathbf{O}_{n \times m} & (\mathbf{c} - \mathbf{A}^\top \boldsymbol{\lambda}_0) \end{bmatrix}, \quad \tilde{\mathbf{b}} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \\ \mathbf{c} \end{bmatrix}$$

(the subscripts above refer to the dimensions/sizes of the corresponding matrices/vectors). Observe that the following point is a strictly interior feasible point for the problem above:

$$\begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \\ \mathbf{u} \\ \mathbf{v} \\ z \end{bmatrix} = \begin{bmatrix} \mathbf{x}_0 \\ \boldsymbol{\lambda}_0 \\ \mathbf{u}_0 \\ \mathbf{v}_0 \\ 1 \end{bmatrix}.$$

Furthermore, the minimum value of the objective function for Karmarkar's artificial problem is zero if and only if the previous set of relations has a solution, that is, there exists  $\mathbf{x}$ ,  $\boldsymbol{\lambda}$ ,  $\mathbf{u}$ , and  $\mathbf{v}$  satisfying

$$\begin{aligned} \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \boldsymbol{\lambda} &= 0, \\ \mathbf{A}\mathbf{x} - \mathbf{v} &= \mathbf{b}, \\ \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{u} &= \mathbf{c}, \\ \mathbf{x}, \boldsymbol{\lambda}, \mathbf{u}, \mathbf{v} &\geq \mathbf{0}. \end{aligned}$$

Therefore, Karmarkar's artificial LP problem is equivalent to the original LP problem:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Note that the main difference between the original LP problem and Karmarkar's artificial problem is that we have an explicit strictly interior feasible point for Karmarkar's artificial problem, and hence we have satisfied the assumption that we imposed at the beginning of this subsection.

## The Algorithm

We are now ready to describe Karmarkar's algorithm. Keep in mind that the LP problem we are solving is a Karmarkar's restricted problem, that is, a problem in Karmarkar's canonical form and satisfies assumptions A, B, C, and D. For convenience, we restate the problem:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n \\ & \text{subject to} && \mathbf{x} \in \Omega \cap \Delta, \end{aligned}$$

where  $\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}\}$  and  $\Delta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$ . Karmarkar's algorithm is an iterative algorithm that, given an initial point  $\mathbf{x}^{(0)}$  and parameter  $q$ , generates a sequence  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$ . Karmarkar's algorithm is described by the following steps:

1. **Initialize:** Set  $k := 0$ ;  $\mathbf{x}^{(0)} = \mathbf{a}_0 = \mathbf{e}/n$ .
2. **Update:** Set  $\mathbf{x}^{(k+1)} = \Psi(\mathbf{x}^{(k)})$ , where  $\Psi$  is an update map described below.
3. **Check the stopping criterion:** If the condition  $\mathbf{c}^\top \mathbf{x}^{(k)}/\mathbf{c}^\top \mathbf{x}^{(0)} \leq 2^{-q}$  is satisfied, then stop.
4. **Iterate:** Set  $k := k + 1$ ; go to step 2.

We describe the update map  $\Psi$  as follows. First, consider the first step in the algorithm:  $\mathbf{x}^{(0)} = \mathbf{a}_0$ . To compute  $\mathbf{x}^{(1)}$ , we use the familiar update equation

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)},$$

where  $\alpha$  is a step size and  $\mathbf{d}^{(0)}$  is an update direction. The step size  $\alpha$  is chosen to be a value in  $(0, 1)$ . Karmarkar recommends a value of  $1/4$  in his original paper [71]. The update direction  $\mathbf{d}^{(0)}$  is chosen as follows. First, note that the gradient of the objective function is  $\mathbf{c}$ . Therefore, the direction of maximum rate of decrease of the objective function is  $-\mathbf{c}$ . However, in general, we cannot simply update along this direction, since  $\mathbf{x}^{(1)}$  is required to lie in the constraint set

$$\begin{aligned} \Omega \cap \Delta &= \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = \mathbf{0}, \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^n : \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \mathbf{x} \geq \mathbf{0} \right\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^n : \mathbf{B}_0 \mathbf{x} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \mathbf{x} \geq \mathbf{0} \right\}, \end{aligned}$$

where  $\mathbf{B}_0 \in \mathbb{R}^{(m+1) \times n}$  is given by

$$\mathbf{B}_0 = \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix}.$$

Note that since  $\mathbf{x}^{(0)} \in \Omega \cap \Delta$ , then for  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}$  also to lie in  $\Omega \cap \Delta$ , the vector  $\mathbf{d}^{(0)}$  must be an element of the nullspace of  $\mathbf{B}_0$ . Hence, we choose  $\mathbf{d}^{(0)}$  to be in the direction of the orthogonal projection of  $-\mathbf{c}$  onto the nullspace of  $\mathbf{B}_0$ . This projection is accomplished by the matrix  $\mathbf{P}_0$  given by

$$\mathbf{P}_0 = \mathbf{I}_n - \mathbf{B}_0^\top (\mathbf{B}_0 \mathbf{B}_0^\top)^{-1} \mathbf{B}_0.$$

Note that  $\mathbf{B}_0 \mathbf{B}_0^\top$  is nonsingular by assumption C. Specifically, we choose  $\mathbf{d}^{(0)}$  to be the vector  $\mathbf{d}^{(0)} = -r \hat{\mathbf{c}}^{(0)}$ , where

$$\hat{\mathbf{c}}^{(0)} = \frac{\mathbf{P}_0 \mathbf{c}}{\|\mathbf{P}_0 \mathbf{c}\|}$$

and  $r = 1/\sqrt{n(n-1)}$ . The scalar  $r$  is incorporated into the update vector  $\mathbf{d}^{(0)}$  for the following reason. First, observe that  $r$  is the radius of the largest sphere inscribed in the simplex  $\Delta$  (see Exercise 18.7). Therefore, the vector  $\mathbf{d}^{(0)} = r \hat{\mathbf{c}}^{(0)}$  points in the direction of the projection  $\hat{\mathbf{c}}^{(0)}$  of  $\mathbf{c}$  onto the nullspace of  $\mathbf{B}_0$  and  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}$  is guaranteed to lie in the constraint set  $\Omega \cap \Delta$ . In fact,  $\mathbf{x}^{(1)}$  lies in the set  $\Omega \cap \Delta \cap \{\mathbf{x} : \|\mathbf{x} - \mathbf{a}_0\| \leq r\}$ . Finally, we note that  $\mathbf{x}^{(1)}$  is a strictly interior point of  $\Delta$ .

The general update step  $\mathbf{x}^{(k+1)} = \Psi(\mathbf{x}^{(k)})$  is performed as follows. We first give a brief description of the basic idea, which is similar to the update from  $\mathbf{x}^{(0)}$  to  $\mathbf{x}^{(1)}$  described above. However, note that  $\mathbf{x}^{(k)}$  is, in general, not at the center of the simplex. Therefore, let us first transform this point to the center. To do this, let  $\mathbf{D}_k$  be a diagonal matrix whose diagonal entries are the components of the vector  $\mathbf{x}^{(k)}$ :

$$\mathbf{D}_k = \begin{bmatrix} x_1^{(k)} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & x_n^{(k)} \end{bmatrix}.$$

It turns out that because  $\mathbf{x}^{(0)}$  is a strictly interior point of  $\Delta$ ,  $\mathbf{x}^{(k)}$  is a strictly interior point of  $\Delta$  for all  $k$  (see Exercise 18.10). Therefore,  $\mathbf{D}_k$  is nonsingular and

$$\mathbf{D}_k^{-1} = \begin{bmatrix} 1/x_1^{(k)} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1/x_n^{(k)} \end{bmatrix}.$$

Consider the mapping  $\mathbf{U}_k : \Delta \rightarrow \Delta$  given by  $\mathbf{U}_k(\mathbf{x}) = \mathbf{D}_k^{-1} \mathbf{x} / \mathbf{e}^\top \mathbf{D}_k^{-1} \mathbf{x}$ . Note that  $\mathbf{U}_k(\mathbf{x}^{(k)}) = \mathbf{e}/n = \mathbf{a}_0$ . We use  $\mathbf{U}_k$  to change the variable from

$\mathbf{x}$  to  $\bar{\mathbf{x}} = \mathbf{U}_k(\mathbf{x})$ . We do this so that  $\mathbf{x}^{(k)}$  is mapped into the center of the simplex, as indicated above. Note that  $\mathbf{U}_k$  is an invertible mapping, with  $\mathbf{x} = \mathbf{U}_k^{-1}(\bar{\mathbf{x}}) = \mathbf{D}_k\bar{\mathbf{x}}/\mathbf{e}^\top \mathbf{D}_k\bar{\mathbf{x}}$ . Letting  $\bar{\mathbf{x}}^{(k)} = \mathbf{U}_k(\mathbf{x}^{(k)}) = \mathbf{a}_0$ , we can now apply the procedure that we described before for getting  $\mathbf{x}^{(1)}$  from  $\mathbf{x}^{(0)} = \mathbf{a}_0$ . Specifically, we update  $\bar{\mathbf{x}}^{(k)}$  to obtain  $\bar{\mathbf{x}}^{(k+1)}$  using the update formula  $\bar{\mathbf{x}}^{(k+1)} = \bar{\mathbf{x}}^{(k)} + \alpha \mathbf{d}^{(k)}$ . To compute  $\mathbf{d}^{(k)}$ , we need to state the original LP problem in the new variable  $\bar{\mathbf{x}}$ :

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{D}_k \bar{\mathbf{x}} \\ & \text{subject to} && \mathbf{A}\mathbf{D}_k \bar{\mathbf{x}} = \mathbf{0} \\ & && \bar{\mathbf{x}} \in \Delta. \end{aligned}$$

The reader can easily verify that the LP problem above in the new variable  $\bar{\mathbf{x}}$  is equivalent to the original LP problem in the sense that  $\mathbf{x}^*$  is an optimal solution to the original problem if and only if  $\mathbf{U}_k(\mathbf{x}^*)$  is an optimal solution to the transformed problem. To see this, simply note that  $\bar{\mathbf{x}} = \mathbf{U}_k(\mathbf{x}) = \mathbf{D}_k^{-1}\mathbf{x}/\mathbf{e}^\top \mathbf{D}_k^{-1}\mathbf{x}$ , and rewrite the objective function and constraints accordingly (see Exercise 18.8). As before, let

$$\mathbf{B}_k = \begin{bmatrix} \mathbf{A}\mathbf{D}_k \\ \mathbf{e}^\top \end{bmatrix}.$$

We choose  $\mathbf{d}^{(k)} = -r\hat{\mathbf{c}}^{(k)}$ , where  $\hat{\mathbf{c}}^{(k)}$  is the normalized projection of  $-(\mathbf{c}^\top \mathbf{D}_k)^\top = -\mathbf{D}_k \mathbf{c}$  onto the nullspace of  $\mathbf{B}_k$ , and  $r = 1/\sqrt{n(n-1)}$  as before. To determine  $\hat{\mathbf{c}}^{(k)}$ , we define the projector matrix  $\mathbf{P}_k$  by

$$\mathbf{P}_k = \mathbf{I}_n - \mathbf{B}_k^\top (\mathbf{B}_k \mathbf{B}_k^\top)^{-1} \mathbf{B}_k.$$

Note that  $\mathbf{B}_k \mathbf{B}_k^\top$  is nonsingular (see Exercise 18.9). The vector  $\hat{\mathbf{c}}^{(k)}$  is therefore given by

$$\hat{\mathbf{c}}^{(k)} = \frac{\mathbf{P}_k \mathbf{D}_k \mathbf{c}}{\|\mathbf{P}_k \mathbf{D}_k \mathbf{c}\|}.$$

The direction vector  $\mathbf{d}^{(k)}$  is then

$$\mathbf{d}^{(k)} = -r\hat{\mathbf{c}}^{(k)} = -r \frac{\mathbf{P}_k \mathbf{D}_k \mathbf{c}}{\|\mathbf{P}_k \mathbf{D}_k \mathbf{c}\|}.$$

The updated vector  $\bar{\mathbf{x}}^{(k+1)} = \bar{\mathbf{x}}^{(k)} + \alpha \mathbf{d}^{(k)}$  is guaranteed to lie in the transformed feasible set  $\{\bar{\mathbf{x}} : \mathbf{A}\mathbf{D}_k \bar{\mathbf{x}} = \mathbf{0}\} \cap \Delta$ . The final step is to apply the inverse transformation  $\mathbf{U}_k^{-1}$  to obtain  $\mathbf{x}^{(k+1)}$ :

$$\mathbf{x}^{(k+1)} = \mathbf{U}_k^{-1}(\bar{\mathbf{x}}^{(k+1)}) = \frac{\mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}{\mathbf{e}^\top \mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}.$$

Note that  $\mathbf{x}^{(k+1)}$  lies in the set  $\Omega \cap \Delta$ . Indeed, we have already seen that  $\mathbf{U}_k$  and  $\mathbf{U}_k^{-1}$  map  $\Delta$  into  $\Delta$ . To see that  $\mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{0}$ , we simply premultiply the foregoing expression by  $\mathbf{A}$  and use the fact that  $\mathbf{A}\mathbf{D}_k \bar{\mathbf{x}}^{(k+1)} = \mathbf{0}$ .

We now summarize the update  $\mathbf{x}^{(k+1)} = \Psi(\mathbf{x}^{(k)})$ :

1. Compute the matrices:

$$\mathbf{D}_k = \begin{bmatrix} x_1^{(k)} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & x_n^{(k)} \end{bmatrix},$$

$$\mathbf{B}_k = \begin{bmatrix} \mathbf{A}\mathbf{D}_k \\ \mathbf{e}^\top \end{bmatrix}.$$

2. Compute the orthogonal projector onto the nullspace of  $\mathbf{B}_k$ :

$$\mathbf{P}_k = \mathbf{I}_n - \mathbf{B}_k^\top (\mathbf{B}_k \mathbf{B}_k^\top)^{-1} \mathbf{B}_k.$$

3. Compute the normalized orthogonal projection of  $\mathbf{c}$  onto the nullspace of  $\mathbf{B}_k$ :

$$\hat{\mathbf{c}}^{(k)} = \frac{\mathbf{P}_k \mathbf{D}_k \mathbf{c}}{\|\mathbf{P}_k \mathbf{D}_k \mathbf{c}\|}.$$

4. Compute the direction vector:

$$\mathbf{d}^{(k)} = -r\hat{\mathbf{c}}^{(k)},$$

where  $r = 1/\sqrt{n(n-1)}$ .

5. Compute  $\bar{\mathbf{x}}^{(k+1)}$  using

$$\bar{\mathbf{x}}^{(k+1)} = \mathbf{a}_0 + \alpha \mathbf{d}^{(k)},$$

where  $\alpha$  is the prespecified step size,  $\alpha \in (0, 1)$ .

6. Compute  $\mathbf{x}^{(k+1)}$  by applying the inverse transformation  $\mathbf{U}_k^{-1}$ :

$$\mathbf{x}^{(k+1)} = \mathbf{U}_k^{-1}(\bar{\mathbf{x}}^{(k+1)}) = \frac{\mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}{\mathbf{e}^\top \mathbf{D}_k \bar{\mathbf{x}}^{(k+1)}}.$$

The matrix  $\mathbf{P}_k$  in step 2 is needed solely for computing  $\mathbf{P}_k \mathbf{D}_k \mathbf{c}$  in step 3. In fact, the two steps can be combined in an efficient way without having to compute  $\mathbf{P}_k$  explicitly, as follows. We first solve a set of linear equations  $\mathbf{B}_k \mathbf{B}_k^\top \mathbf{y} = \mathbf{B}_k \mathbf{D}_k \mathbf{c}$  (for the variable  $\mathbf{y}$ ), and then compute  $\mathbf{P}_k \mathbf{D}_k \mathbf{c}$  using the expression  $\mathbf{P}_k \mathbf{D}_k \mathbf{c} = \mathbf{D}_k \mathbf{c} - \mathbf{B}_k^\top \mathbf{y}$ .

For more details on Karmarkar's algorithm, see [42], [55], [71], and [124]. For an informal introduction to the algorithm, see [110]. For further reading on other nonsimplex methods in linear programming, see [42], [55], [96], and [119]. A continuous gradient system for solving linear programming problems is discussed in [26]. An interesting three-article series on developments of

the linear programming area before and after 1984 appeared in *SIAM News*, Vol. 22, No. 2, March 1989. The first article in this journal issue contains an account by Wright on recent progress and a history of linear programming from the early 1800s. The second article, by Anstreicher, focuses on interior-point algorithms developed since 1984. Finally in the third article in the series, Monma surveys computational implementations of interior-point methods.

## EXERCISES

**18.1** Write a simple MATLAB function to implement the affine scaling algorithm. The inputs are  $\mathbf{c}$ ,  $\mathbf{A}$ ,  $\mathbf{b}$ , and  $\mathbf{x}^{(0)}$ , where  $\mathbf{x}^{(0)}$  is a strictly feasible initial point. Test the function on the problem in Example 16.2; use  $\mathbf{x}^{(0)} = [2, 3, 2, 3, 3]^\top$ .

**18.2** Write a MATLAB routine that implements the two-phase affine scaling method. It may be useful to use the MATLAB function of Exercise 18.1. Test the routine on the problem in Example 16.5.

**18.3** For a given linear programming problem of the form

$$\begin{aligned} &\text{minimize} \quad \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} \quad \mathbf{Ax} \geq \mathbf{b} \\ &\quad \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

the associated Karmarkar's artificial problem can be solved directly using the affine scaling method. Write a simple MATLAB program to solve problems of the form above by using the affine scaling algorithm applied to the associated Karmarkar's artificial problem. It may be useful to use the MATLAB function of Exercise 18.1. Test your program on the problem in Example 15.15.

**18.4** Let  $\mathbf{a} \in \mathbb{R}^n$ ,  $\mathbf{a} > \mathbf{0}$ . Let  $\mathbf{T} = [T_1, \dots, T_{n+1}]$  be the projective transformation of the positive orthant  $P_+$  of  $\mathbb{R}^n$  into the simplex  $\Delta$  in  $\mathbb{R}^{n+1}$ , given by

$$T_i(\mathbf{x}) = \begin{cases} \frac{x_i/a_i}{x_1/a_1 + \dots + x_n/a_n + 1} & \text{if } 1 \leq i \leq n \\ \frac{1}{x_1/a_1 + \dots + x_n/a_n + 1} & \text{if } i = n+1. \end{cases}$$

Prove the following properties of  $\mathbf{T}$  (see [71]):

1.  $\mathbf{T}$  is a one-to-one mapping; that is,  $\mathbf{T}(\mathbf{x}) = \mathbf{T}(\mathbf{y})$  implies that  $\mathbf{x} = \mathbf{y}$ .
2.  $\mathbf{T}$  maps  $P_+$  onto  $\Delta \setminus \{\mathbf{x} : x_{n+1} = 0\} \triangleq \{\mathbf{x} \in \Delta : x_{n+1} > 0\}$ ; that is, for each  $\mathbf{y} \in \{\mathbf{x} \in \Delta : x_{n+1} > 0\}$ , there exists  $\mathbf{x} \in P_+$  such that  $\mathbf{y} = \mathbf{T}(\mathbf{x})$ .
3. The inverse transformation of  $\mathbf{T}$  exists on  $\{\mathbf{x} \in \Delta : x_{n+1} > 0\}$  and is given by  $\mathbf{T}^{-1} = [T_1^{-1}, \dots, T_n^{-1}]^\top$ , with  $T_i^{-1}(\mathbf{y}) = a_i y_i / y_{n+1}$ .

4.  $\mathbf{T}$  maps  $\mathbf{a}$  to the center of the simplex  $\Delta$ , that is,  $\mathbf{T}(\mathbf{a}) = \mathbf{e}/(n+1) = [1/(n+1), \dots, 1/(n+1)] \in \mathbb{R}^{n+1}$ .
5. Suppose that  $\mathbf{x}$  satisfies  $\mathbf{Ax} = \mathbf{b}$ , and  $\mathbf{y} = \mathbf{T}(\mathbf{x})$ . Let  $\mathbf{x}' = [y_1 a_1, \dots, y_n a_n]^\top$ . Then,  $\mathbf{Ax}' = \mathbf{b} y_{n+1}$ .

**18.5** Let  $\mathbf{T}$  be the projective transformation in Exercise 18.4 and  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a given matrix. Prove that there exists a matrix  $\mathbf{A}' \in \mathbb{R}^{m \times (n+1)}$  such that  $\mathbf{Ax} = \mathbf{b}$  if and only if  $\mathbf{A}'\mathbf{T}(\mathbf{x}) = \mathbf{0}$ .

*Hint:* Let the  $i$ th column of  $\mathbf{A}'$  be given by  $a_i$  times the  $i$ th column of  $\mathbf{A}$ ,  $i = 1, \dots, n$ , and the  $(n+1)$ th column of  $\mathbf{A}'$  be given by  $-\mathbf{b}$ .

**18.6** Let  $\mathbf{T}$  be the projective transformation in Exercise 18.4 and  $\mathbf{c} \in \mathbb{R}^n$  be a given vector. Prove that there exists a vector  $\mathbf{c}' \in \mathbb{R}^{n+1}$  such that  $\mathbf{c}^\top \mathbf{x} = 0$  if and only if  $\mathbf{c}'^\top \mathbf{T}(\mathbf{x}) = 0$ .

*Hint:* Use property 3 in Exercise 18.4, with the  $\mathbf{c}' = [c'_1, \dots, c'_{n+1}]^\top$  given by  $c'_i = a_i c_i$ ,  $i = 1, \dots, n$ , and  $c'_{n+1} = 0$ .

**18.7** Let  $\Delta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^\top \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\}$  be the simplex in  $\mathbb{R}^n$ ,  $n > 1$ , and let  $\mathbf{a}_0 = \mathbf{e}/n$  be its center. A *sphere* of radius  $r$  centered at  $\mathbf{a}_0$  is the set  $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{a}_0\| \leq r\}$ . The sphere is said to be *inscribed in*  $\Delta$  if  $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{a}_0\| = r, \mathbf{e}^\top \mathbf{x} = 1\} \subset \Delta$ . Show that the largest such sphere has radius  $r = 1/\sqrt{n(n-1)}$ .

**18.8** Consider the following Karmarkar's restricted problem:

$$\begin{aligned} &\text{minimize} && \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} && \mathbf{Ax} = \mathbf{0} \\ &&& \mathbf{x} \in \Delta. \end{aligned}$$

Let  $\mathbf{x}_0 \in \Delta$  be a strictly interior point of  $\Delta$ , and  $\mathbf{D}$  be a diagonal matrix whose diagonal entries are the components of  $\mathbf{x}_0$ . Define the map  $\mathbf{U} : \Delta \rightarrow \Delta$  by  $\mathbf{U}(\mathbf{x}) = \mathbf{D}^{-1}\mathbf{x}/\mathbf{e}^\top \mathbf{D}^{-1}\mathbf{x}$ . Let  $\bar{\mathbf{x}} = \mathbf{U}(\mathbf{x})$  represent a change of variable. Show that the following transformed LP problem in the variable  $\bar{\mathbf{x}}$ ,

$$\begin{aligned} &\text{minimize} && \mathbf{c}^\top \mathbf{D}\bar{\mathbf{x}} \\ &\text{subject to} && \mathbf{AD}\bar{\mathbf{x}} = \mathbf{0} \\ &&& \bar{\mathbf{x}} \in \Delta, \end{aligned}$$

is equivalent to the original LP problem above in the sense that  $\mathbf{x}^*$  is an optimal solution to the original problem if and only if  $\bar{\mathbf{x}}^* = \mathbf{U}(\mathbf{x}^*)$  is an optimal solution to the transformed problem.

**18.9** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ , and  $\Omega = \{\mathbf{x} : \mathbf{Ax} = \mathbf{0}\}$ . Suppose that  $\mathbf{A}$  satisfies

$$\text{rank} \begin{bmatrix} \mathbf{A} \\ \mathbf{e}^\top \end{bmatrix} = m + 1.$$

Let  $\mathbf{x}_0 \in \Delta \cap \Omega$  be a strictly interior point of  $\Delta \subset \mathbb{R}^n$  and  $\mathbf{D}$  be a diagonal matrix whose diagonal entries are the components of  $\mathbf{x}_0$ . Consider the matrix  $\mathbf{B}$  defined by

$$\mathbf{B} = \begin{bmatrix} \mathbf{AD} \\ \mathbf{e}^\top \end{bmatrix}.$$

Show that  $\text{rank } \mathbf{B} = m + 1$ , and hence  $\mathbf{BB}^\top$  is nonsingular.

**18.10** Show that in Karmarkar's algorithm,  $\mathbf{x}^{(k)}$  is a strictly interior point of  $\Delta$ .

# CHAPTER 19

---

## INTEGER LINEAR PROGRAMMING

---

### 19.1 Introduction

This chapter is devoted to linear programs with the additional constraint that the solution components be integers. Such problems are called *integer linear programming (ILP)* (or simply *integer programming*) problems, and arise naturally in many practical situations. For example, in Example 15.1, the decision variables represent production levels, which we allowed to take real values. If production levels correspond to actual numbers of products, then it is natural to impose the constraint that they be integer valued. If we expect solutions that are very large in magnitude, then ignoring the integer constraint might have little practical consequence. However, in cases where the solution is a relatively small integer (on the order of 10, say), then ignoring the integer constraint could lead to dramatically erroneous solutions.

Throughout this section, we use the notation  $\mathbb{Z}$  for the set of integers,  $\mathbb{Z}^n$  the set of vectors with  $n$  integer components, and  $\mathbb{Z}^{m \times n}$  the set of  $m \times n$  matrices with integer entries. Using this notation, we can express an ILP

problem in following form:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \\ & && \mathbf{x} \in \mathbb{Z}^n. \end{aligned}$$

## 19.2 Unimodular Matrices

There is a class of ILP problems that can be solved using standard linear programming methods. To proceed, we need some definitions and background results. The reader should recall the definition of a *minor* from Section 2.2.

**Definition 19.1** An  $m \times n$  integer matrix  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ , is *unimodular* if all its nonzero  $m$ th-order minors are  $\pm 1$  (i.e., either 1 or  $-1$ ). ■

Unimodular matrices play a special role in the context of linear equations and integer basic solutions. Consider the linear equation  $\mathbf{Ax} = \mathbf{b}$  with  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ . Let  $\mathbf{B}$  be a corresponding basis matrix (an  $m \times m$  matrix consisting of  $m$  linearly independent columns of  $\mathbf{A}$ ). Then, the unimodularity of  $\mathbf{A}$  is equivalent to  $|\det \mathbf{B}| = 1$  for any such  $\mathbf{B}$ . The following lemma connects unimodularity with integer basic solutions.

**Lemma 19.1** Consider the linear equation  $\mathbf{Ax} = \mathbf{b}$  where  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ , is unimodular and  $\mathbf{b} \in \mathbb{Z}^m$ . Then, all basic solutions have integer components. □

*Proof.* As usual, suppose that the first  $m$  columns of  $\mathbf{A}$  constitute a basis, and that  $\mathbf{B}$  is the invertible  $m \times m$  matrix composed of these columns. Then the corresponding basic solution is

$$\mathbf{x}^* = \begin{bmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0} \end{bmatrix}.$$

Because all the elements of  $\mathbf{A}$  are integers,  $\mathbf{B}$  is an integer matrix. Moreover, because  $\mathbf{A}$  is unimodular,  $|\det \mathbf{B}| = 1$ . This implies that the inverse  $\mathbf{B}^{-1}$  is also an integer matrix (see [62, p. 21]). Therefore,  $\mathbf{x}^*$  is an integer vector. ■

**Corollary 19.1** Consider the LP constraint

$$\begin{aligned} & \mathbf{Ax} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{A}$  is unimodular,  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ , and  $\mathbf{b} \in \mathbb{Z}^m$ . Then, all basic feasible solutions have integer components. □

Unimodularity allows us to solve ILP problems using the simplex method. Specifically, consider the ILP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \\ & && \mathbf{x} \in \mathbb{Z}^n \end{aligned}$$

where  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ , is unimodular and  $\mathbf{b} \in \mathbb{Z}^m$ . Then, the corollary above tells us that if we consider the associated LP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

the optimal basic feasible solution is an integer vector. This means that we can apply the simplex method to the LP problem above to obtain a solution to the original ILP problem.

**Example 19.1** Consider the following ILP problem:

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 \\ & \text{subject to} && x_1 + x_3 = 4 \\ & && x_2 + x_4 = 6 \\ & && x_1 + x_2 + x_5 = 8 \\ & && x_1, x_2, x_3, x_4, x_5 \geq 0 \\ & && x_1, x_2, x_3, x_4, x_5 \in \mathbb{Z} \end{aligned}$$

We can write this problem in matrix form with

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 4 \\ 6 \\ 8 \end{bmatrix}.$$

Notice that  $\mathbf{b} \in \mathbb{Z}^3$ . Moreover, it is easy to check that  $\mathbf{A}$  is unimodular. Hence, the ILP problem above can be solved by solving the LP problem

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 \\ & \text{subject to} && x_1 + x_3 = 4 \\ & && x_2 + x_4 = 6 \\ & && x_1 + x_2 + x_5 = 8 \\ & && x_1, x_2, x_3, x_4, x_5 \geq 0. \end{aligned}$$

This was done in Example 16.2 using the simplex method, yielding optimal solution  $[2, 6, 2, 0, 0]^\top$ , which is an integer vector. ■

In general, when the matrix  $\mathbf{A}$  is not unimodular, the simplex method applied to the associated LP problem yields a noninteger optimal solution. However, in some cases, even if  $\mathbf{A}$  is not unimodular, the simplex method still produces an integer optimal basic feasible solution. To see this, suppose that we are given  $\mathbf{A} \in \mathbb{Z}^{m \times n}$ ,  $m \leq n$ , and  $\mathbf{b} \in \mathbb{Z}^m$ . Note that as long as each  $m \times m$  basis matrix  $\mathbf{B}$  consisting of columns of  $\mathbf{A}$  corresponding to a basic *feasible* solution has the property that  $|\det \mathbf{B}| = 1$ , we can use the argument in the proof of Lemma 19.1 to conclude that the basic feasible solution is an integer vector. Equivalently, we can draw this conclusion if each basis submatrix  $\mathbf{B}$  of  $\mathbf{A}$  such that  $|\det \mathbf{B}| \neq 1$  corresponds to a *nonfeasible* basic solution. We illustrate this in the following example.

**Example 19.2** Consider the ILP problem

$$\begin{aligned} & \text{minimize} && -x_1 - 2x_2 \\ & \text{subject to} && -2x_1 + x_2 + x_3 = 2 \\ & && -x_1 + x_2 + x_4 = 3 \\ & && x_1 + x_5 = 3 \\ & && x_i \geq 0, \quad i = 1, \dots, 5 \\ & && x_i \in \mathbb{Z}, \quad i = 1, \dots, 5. \end{aligned}$$

Can this ILP problem be solved using the simplex method? We can easily verify that the matrix

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

is not unimodular. Indeed, it has one (and only one) basis submatrix with determinant other than  $\pm 1$ , consisting of the first, fourth, and fifth columns of  $\mathbf{A}$ . Indeed, if we write  $\mathbf{B} = [\mathbf{a}_1, \mathbf{a}_4, \mathbf{a}_5]$ , then  $\det \mathbf{B} = -2$ . However, a closer examination of this matrix and the vector  $\mathbf{b} = [2, 3, 3]^\top$  reveals that the corresponding basic solution is not feasible:  $\mathbf{B}^{-1}\mathbf{b} = [-1, 2, 4]^\top$  (which, coincidentally, happens to be an integer vector). Therefore, for this problem, applying the simplex method to the associated LP problem will produce an integer optimal basic feasible solution, which also solves the original ILP problem.

We begin by forming the first tableau,

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ -2 & 1 & 1 & 0 & 0 & 2 \\ -1 & 1 & 0 & 1 & 0 & 3 \\ 1 & 0 & 0 & 0 & 1 & 3 \\ \mathbf{c}^\top & -1 & -2 & 0 & 0 & 0 & 0 \end{array}$$

We have  $r_2 = -2$ . Therefore, we introduce  $\mathbf{a}_2$  into the new basis. We calculate the ratios  $y_{i0}/y_{i2}$ ,  $y_{i2} > 0$ , to determine the pivot element:

$$\frac{y_{10}}{y_{12}} = \frac{2}{1} \quad \text{and} \quad \frac{y_{20}}{y_{22}} = \frac{3}{1}.$$

We will use  $y_{12}$  as the pivot. Performing elementary row operations, we obtain the second tableau,

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ -2 & 1 & 1 & 0 & 0 & 2 \\ 1 & 0 & -1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 3 \\ \mathbf{r}^\top & -5 & 0 & 2 & 0 & 0 & 4 \end{array}$$

We now have  $r_1 = -5 < 0$ . Therefore, we introduce  $\mathbf{a}_1$  into the new basis. We next calculate the ratios  $y_{i0}/y_{i2}$ ,  $y_{i2} > 0$ , to determine the pivot element:

$$\frac{y_{20}}{y_{21}} = \frac{1}{1} \quad \text{and} \quad \frac{y_{30}}{y_{31}} = \frac{3}{1}.$$

We will use  $y_{21}$  as the pivot. Performing row elementary operations, we obtain the third tableau,

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 0 & 1 & -1 & 2 & 0 & 4 \\ 1 & 0 & -1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 & 2 \\ \mathbf{r}^\top & 0 & 0 & -3 & 5 & 0 & 9 \end{array}$$

We have  $r_3 = -3 < 0$ . Therefore, we introduce  $\mathbf{a}_3$  into the new basis. We next calculate the ratios  $y_{i0}/y_{i2}$ ,  $y_{i2} > 0$ , to determine the pivot element,

$$\frac{y_{30}}{y_{33}} = \frac{2}{1}.$$

We will use  $y_{33}$  as the pivot. Performing row elementary operations, we obtain the fourth tableau,

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 0 & 1 & 0 & 1 & 1 & 6 \\ 1 & 0 & 0 & 0 & 1 & 3 \\ 0 & 0 & 1 & -1 & 1 & 2 \\ \mathbf{r}^\top & 0 & 0 & 0 & 2 & 3 & 15 \end{array}$$

All reduced cost coefficients are now positive, which means that the current solution is  $[3, 6, 2, 0, 0]^\top$ . ■

Next, we consider ILP problems of the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \\ & && \mathbf{x} \in \mathbb{Z}^n \end{aligned}$$

We have seen in Section 15.5 that we can transform the inequality constraint  $\mathbf{A}\mathbf{x} \leq \mathbf{b}$  into standard form by introducing slack variables. Doing so would lead to a new problem in standard form for which the constraint has the form  $[\mathbf{A}, \mathbf{I}]\mathbf{y} = \mathbf{b}$  (where the vector  $\mathbf{y}$  contains  $\mathbf{x}$  and the slack variables). To deal with matrices of the form  $[\mathbf{A}, \mathbf{I}]$ , we need another definition.

**Definition 19.2** An  $m \times n$  integer matrix  $\mathbf{A} \in \mathbb{Z}^{m \times n}$  is *totally unimodular* if all its nonzero minors are  $\pm 1$ . ■

By minors here we mean  $p$ th-order minors for  $p \leq \min(m, n)$ . Equivalently, a matrix  $\mathbf{A} \in \mathbb{Z}^{m \times n}$  is totally unimodular if and only if all its square invertible submatrices have determinant  $\pm 1$ . By a *submatrix* of  $\mathbf{A}$  we mean a matrix obtained by removing some columns and rows of  $\mathbf{A}$ . It is easy to see from this definition that if an integer matrix is totally unimodular, then each entry is 0, 1, or  $-1$ . The next proposition relates the total unimodularity of  $\mathbf{A}$  with the unimodularity of  $[\mathbf{A}, \mathbf{I}]$  (see also Exercise 19.3).

**Proposition 19.1** If an  $m \times n$  integer matrix  $\mathbf{A} \in \mathbb{Z}^{m \times n}$  is totally unimodular, then the matrix  $[\mathbf{A}, \mathbf{I}]$  is unimodular. □

*Proof.* Let  $\mathbf{A}$  satisfy the assumptions of the proposition. We will show that any  $m \times m$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  has determinant  $\pm 1$ . We first note that any  $m \times m$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  that consists only of columns of  $\mathbf{A}$  has determinant  $\pm 1$  because  $\mathbf{A}$  is totally unimodular. Moreover, the  $m \times m$  submatrix  $\mathbf{I}$  satisfies  $\det \mathbf{I} = 1$ .

Consider now an  $m \times m$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  composed of  $k$  columns of  $\mathbf{A}$  and  $m - k$  columns of  $\mathbf{I}$ . Without loss of generality, suppose that this submatrix is composed of the last  $k$  columns of  $\mathbf{A}$  and the first

$m - k$  columns of  $\mathbf{I}$ ; that is, the  $m \times m$  invertible submatrix is

$$\mathbf{B} = \begin{bmatrix} \mathbf{a}_{n-k+1} & \cdots & \mathbf{a}_n & \mathbf{e}_1 & \cdots & \mathbf{e}_{m-k} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_{m-k,k} & \mathbf{I}_{m-k} \\ \mathbf{B}_{k,k} & \mathbf{O} \end{bmatrix},$$

where  $\mathbf{e}_i$  is the  $i$ th column of the identity matrix. This choice of columns is without loss of generality because we can exchange rows and columns to arrive at this form, and each exchange only changes the sign of the determinant. Moreover, note that  $\det \mathbf{B} = \pm \det \mathbf{B}_{k,k}$  (see also Exercises 19.4 and 2.4). Thus,  $\mathbf{B}_{k,k}$  is invertible because  $\mathbf{B}$  is invertible. Moreover, because  $\mathbf{B}_{k,k}$  is a submatrix of  $\mathbf{A}$  and  $\mathbf{A}$  is totally unimodular,  $\det \mathbf{B}_{k,k} = \pm 1$ . Hence,  $\det \mathbf{B} = \pm 1$  also. Thus any  $m \times m$  invertible submatrix of  $[\mathbf{A}, \mathbf{I}]$  has determinant  $\pm 1$ , which implies that  $[\mathbf{A}, \mathbf{I}]$  is unimodular. ■

Combining the result above with Lemma 19.1, we obtain the following corollary.

**Corollary 19.2** *Consider the LP constraint*

$$\begin{aligned} [\mathbf{A}, \mathbf{I}]\mathbf{x} &= \mathbf{b} \\ \mathbf{x} &\geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{Z}^{m \times n}$  is totally unimodular and  $\mathbf{b} \in \mathbb{Z}^m$ . Then, all basic feasible solutions have integer components. □

Total unimodularity of  $\mathbf{A}$  allows us to solve ILP problems of the following form using the simplex method:

$$\begin{aligned} &\text{minimize } \mathbf{c}^\top \mathbf{x} \\ \text{subject to } &\mathbf{Ax} \leq \mathbf{b} \\ &\mathbf{x} \geq \mathbf{0} \\ &\mathbf{x} \in \mathbb{Z}^n \end{aligned}$$

where  $\mathbf{b} \in \mathbb{Z}^m$ . Specifically, we first consider the associated LP problem

$$\begin{aligned} &\text{minimize } \mathbf{c}^\top \mathbf{x} \\ \text{subject to } &\mathbf{Ax} \leq \mathbf{b} \\ &\mathbf{x} \geq \mathbf{0}. \end{aligned}$$

If  $\mathbf{A}$  is totally unimodular, then the corollary above tells us that once we convert this problem into standard form by introducing a slack-variable vector  $\mathbf{z}$ ,

$$\begin{aligned} &\text{minimize } \mathbf{c}^\top \mathbf{x} \\ \text{subject to } &[\mathbf{A}, \mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} = \mathbf{b} \\ &\mathbf{x}, \mathbf{z} \geq \mathbf{0}, \end{aligned}$$

the optimal basic feasible solution is an integer vector. This means that we can apply the simplex method to the LP problem above to obtain a solution to the original ILP problem. Note that although we only needed the  $\mathbf{x}$  part of the solution to be integer, the slack-variable vector  $\mathbf{z}$  is automatically integer for any integer  $\mathbf{x}$ , because both  $\mathbf{A}$  and  $\mathbf{b}$  only contain integers (see also Exercise 19.5).

**Example 19.3** Consider the following ILP problem:

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 \\ & \text{subject to} && x_1 \leq 4 \\ & && x_2 \leq 6 \\ & && x_1 + x_2 \leq 8 \\ & && x_1, x_2 \geq 0 \\ & && x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

This problem can be written in the matrix form above with

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 4 \\ 6 \\ 8 \end{bmatrix}.$$

It is easy to check that  $\mathbf{A}$  is totally unimodular. Hence, the ILP problem above can be solved by solving the LP problem

$$\begin{aligned} & \text{maximize} && 2x_1 + 5x_2 \\ & \text{subject to} && x_1 + x_3 = 4 \\ & && x_2 + x_4 = 6 \\ & && x_1 + x_2 + x_5 = 8 \\ & && x_1, x_2, x_3, x_4, x_5 \geq 0, \end{aligned}$$

as was done in Example 16.2. ■

As discussed before, even if  $[\mathbf{A}, \mathbf{I}]$  is not unimodular, the simplex algorithm might still yield a solution to the original ILP. In particular, even if  $\mathbf{A}$  is not totally unimodular, the method above might still work, as illustrated in the following example.

**Example 19.4** Consider the following ILP problem:

$$\begin{aligned} & \text{maximize} && x_1 + 2x_2 \\ & \text{subject to} && -2x_1 + x_2 \leq 2 \\ & && x_1 - x_2 \geq -3 \\ & && x_1 \leq 3 \\ & && x_1 \geq 0, x_2 \geq 0, x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

We first express the given problem in this equivalent form:

$$\begin{aligned} \text{minimize} \quad & -x_1 - 2x_2 \\ \text{subject to} \quad & -2x_1 + x_2 \leq 2 \\ & -x_1 + x_2 \leq 3 \\ & x_1 \leq 3 \\ & x_1 \geq 0, \quad x_2 \geq 0, \quad x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

We next represent the problem above in standard form by introducing slack variables  $x_3$ ,  $x_4$ , and  $x_5$  to obtain

$$\begin{aligned} \text{minimize} \quad & -x_1 - 2x_2 \\ \text{subject to} \quad & -2x_1 + x_2 + x_3 = 2 \\ & -x_1 + x_2 + x_4 = 3 \\ & x_1 + x_5 = 3 \\ & x_i \geq 0, \quad i = 1, \dots, 5. \end{aligned}$$

This problem is now of the form in Example 19.2, where the simplex method was used. Recall that the solution is  $[3, 6, 2, 0, 0]^\top$ . Thus, the solution to the original problem is  $\mathbf{x}^* = [3, 6]^\top$ .

Note that the matrix

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

is not totally unimodular, because it has an entry  $(-2)$  not equal to  $0$ ,  $1$ , or  $-1$ . Indeed, the matrix  $[\mathbf{A}, \mathbf{I}]$  is not unimodular. However, in this case, the simplex method still produced an optimal solution to the ILP, as explained in Example 19.2. ■

### 19.3 The Gomory Cutting-Plane Method

In 1958, Ralph E. Gomory [54] proposed a method where noninteger optimal solutions obtained using the simplex method are successively removed from the feasible set by adding constraints that exclude these noninteger solutions from the feasible set. The additional constraints, referred to as *Gomory cuts*, do not eliminate integer feasible solutions from the feasible set. The process is repeated until the optimal solution is an integer vector.

To describe Gomory cuts, we use the *floor* operator, defined next.

**Definition 19.3** The floor of a real number, denoted  $\lfloor x \rfloor$ , is the integer obtained by rounding  $x$  toward  $-\infty$ . ■

For example,  $\lfloor 3.4 \rfloor = 3$  and  $\lfloor -3.4 \rfloor = -4$ .

Consider the ILP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \\ & && \mathbf{x} \in \mathbb{Z}^n. \end{aligned}$$

We begin by applying the simplex method to obtain an optimal basic feasible solution to the LP problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

As usual, suppose that the first  $m$  columns form the basis for the optimal basic feasible solution. The corresponding canonical augmented matrix is

$$\begin{array}{cccccccccc} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_i & \cdots & \mathbf{a}_m & \mathbf{a}_{m+1} & \cdots & \mathbf{a}_n & \mathbf{y}_0 \\ 1 & 0 & \cdots & 0 & \cdots & 0 & y_{1,m+1} & \cdots & y_{1,n} & y_{10} \\ 0 & 1 & \cdots & 0 & \cdots & 0 & y_{2,m+1} & \cdots & y_{2,n} & y_{20} \\ \vdots & \vdots & & \vdots & & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & \cdots & 0 & y_{i,m+1} & \cdots & y_{i,n} & y_{i0} \\ \vdots & \vdots & & \vdots & & \vdots & & & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 & y_{m,m+1} & \cdots & y_{m,n} & y_{m0} \end{array}$$

Consider the  $i$ th component of the optimal basic feasible solution,  $y_{i0}$ . Suppose that  $y_{i0}$  is not an integer. Note that any feasible vector  $\mathbf{x}$  satisfies the equality constraint (taken from the  $i$ th row)

$$x_i + \sum_{j=m+1}^n y_{ij}x_j = y_{i0}.$$

We use this equation to derive an additional constraint that would eliminate the current optimal noninteger solution from the feasible set without eliminating any integer feasible solution. To see how, consider the inequality constraint

$$x_i + \sum_{j=m+1}^n \lfloor y_{ij} \rfloor x_j \leq y_{i0}.$$

Because  $\lfloor y_{ij} \rfloor \leq y_{ij}$ , any  $\mathbf{x} \geq \mathbf{0}$  that satisfies the first equality constraint above also satisfies this inequality constraint. Thus, any feasible  $\mathbf{x}$  satisfies this inequality constraint. Moreover, for any integer feasible vector  $\mathbf{x}$ , the

left-hand side of the inequality constraint is an integer. Therefore, any integer feasible vector  $\mathbf{x}$  also satisfies

$$x_i + \sum_{j=m+1}^n \lfloor y_{ij} \rfloor x_j \leq \lfloor y_{i0} \rfloor.$$

Subtracting this inequality from the equation above, we deduce that any integer feasible vector satisfies

$$\sum_{j=m+1}^n (y_{ij} - \lfloor y_{ij} \rfloor) x_j \geq y_{i0} - \lfloor y_{i0} \rfloor.$$

Next, notice that the optimal basic feasible solution above does not satisfy this inequality, because the left-hand side for the optimal basic feasible solution is 0, but the right-hand side is a positive number. Therefore, if we impose the additional inequality constraint above to the original LP problem, the new constraint set would be such that the current optimal basic feasible solution is no longer feasible, but yet every *integer* feasible vector remains feasible. This new constraint is called a *Gomory cut*.

To transform the new LP problem into standard form, we introduce the surplus variable  $x_{n+1}$  to obtain the equality constraint

$$\sum_{j=m+1}^n (y_{ij} - \lfloor y_{ij} \rfloor) x_j - x_{n+1} = y_{i0} - \lfloor y_{i0} \rfloor.$$

For convenience, we will also call this equality constraint a *Gomory cut*. By augmenting this equation into  $\mathbf{A}$  and  $\mathbf{b}$ , or canonical versions of them (e.g., in the form of a simplex tableau), we obtain a new LP problem in standard form. We can then solve the new problem using the simplex method and examine the resulting optimal basic feasible solution. If the solution satisfies the integer constraints, then we are done—this vector gives an optimal solution to the original ILP problem by extracting the appropriate components. If the solution does not satisfy the integer constraints, we introduce another Gomory cut and repeat the process. We call this procedure the *Gomory cutting-plane method*.

Note that in applying the Gomory cutting-plane method, we only need to introduce enough cuts to satisfy the integer constraints for the original ILP problem. The additional variables introduced by slack variables or by the Gomory cuts are not constrained to be integers.

In the following two examples, we illustrate how the Gomory cutting-plane method can be implemented by incorporating Gomory cuts directly into the simplex tableau.

**Example 19.5** Consider the following ILP problem<sup>1</sup>:

$$\begin{aligned} & \text{maximize} && 3x_1 + 4x_2 \\ & \text{subject to} && \frac{2}{5}x_1 + x_2 \leq 3 \\ & && \frac{2}{5}x_1 - \frac{2}{5}x_2 \leq 1 \\ & && x_1, x_2 \geq 0 \\ & && x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

We first solve the problem graphically. The constraint set  $\Omega$  for the associated LP problem (without integer constraints) can be found by calculating the extreme points:

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0 & 0 \end{bmatrix}^\top, \quad \mathbf{x}^{(2)} = \begin{bmatrix} \frac{5}{2} & 0 \end{bmatrix}^\top, \quad \mathbf{x}^{(3)} = \begin{bmatrix} 0 & 3 \end{bmatrix}^\top, \quad \mathbf{x}^{(4)} = \begin{bmatrix} \frac{55}{14} & \frac{10}{7} \end{bmatrix}^\top.$$

In Figure 19.1, we show the feasible set  $\Omega$ . In Figure 19.2, we show the feasible set for the ILP problem, which allows us to solve the problem graphically. The solution is obtained by finding the straight line  $f = 3x_1 + 4x_2$  with largest  $f$  that passes through a feasible point with integer components. This can be accomplished by first drawing the line  $f = 3x_1 + 4x_2$  for  $f = 0$  and then gradually increasing the values of  $f$ , which corresponds to sliding across the feasible region until the straight line passes through the “last” integer feasible point yielding the largest value of the objective function. From Figure 19.2, we can see that the optimal solution to the ILP problem is  $[2, 2]^\top$ .

We now solve the problem using the Gomory cutting-plane method. First we represent the associated LP problem in standard form:

$$\begin{aligned} & \text{maximize} && 3x_1 + 4x_2 \\ & \text{subject to} && \frac{2}{5}x_1 + x_2 + x_3 = 3 \\ & && \frac{2}{5}x_1 - \frac{2}{5}x_2 + x_4 = 1 \\ & && x_1, x_2, x_3, x_4 \geq 0. \end{aligned}$$

Note that we only need the first two components of the solution to be integers. We can start the simplex method because we have an obvious basic feasible solution. The first tableau is

$$\begin{array}{ccccc|c} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ \frac{2}{5} & 1 & 1 & 0 & 3 \\ \frac{2}{5} & -\frac{2}{5} & 0 & 1 & 1 \\ \mathbf{c}^\top & -3 & -4 & 0 & 0 \end{array}$$

<sup>1</sup>Thanks to David Schvartzman Cohenca for his solution to this problem.

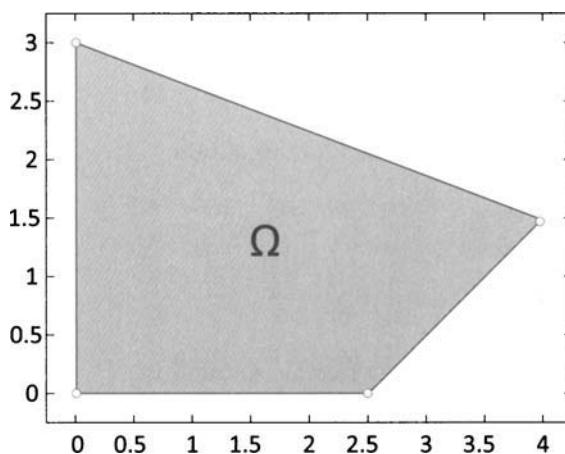


Figure 19.1 Feasible set  $\Omega$  for LP problem in Example 19.5.

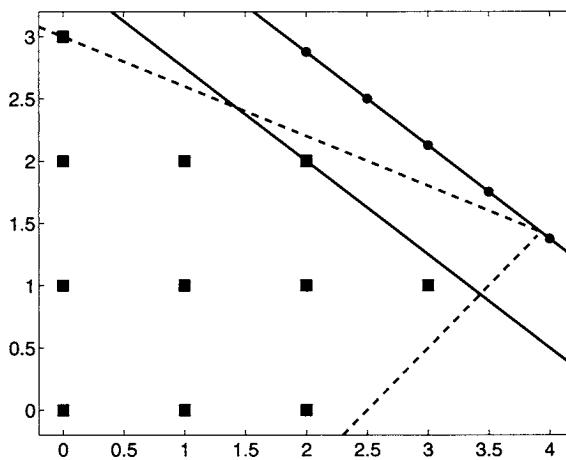


Figure 19.2 Graphical solution for ILP problem in Example 19.5.

We bring  $a_2$  into the basis and pivot about the element (1, 2) to obtain

$$\begin{array}{ccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ \frac{2}{5} & 1 & 1 & 0 & 3 \\ \frac{14}{25} & 0 & \frac{2}{5} & 1 & \frac{11}{5} \\ \mathbf{r}^\top & -\frac{7}{5} & 0 & 4 & 0 & 12 \end{array}$$

Next, we pivot about the element (2, 1) to obtain

$$\begin{array}{ccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ 0 & 1 & \frac{10}{14} & -\frac{10}{14} & \frac{20}{14} \\ 1 & 0 & \frac{10}{14} & \frac{25}{14} & \frac{55}{14} \\ \mathbf{r}^\top & 0 & 0 & 5 & \frac{5}{2} & \frac{35}{2} \end{array}$$

The corresponding optimal basic feasible solution is

$$\left[ \frac{55}{14} \quad \frac{10}{7} \quad 0 \quad 0 \right]^\top,$$

which does not satisfy the integer constraints.

We start by introducing the Gomory cut corresponding to the first row of the tableau. We obtain

$$\frac{10}{14}x_3 + \frac{4}{14}x_4 - x_5 = \frac{6}{14}.$$

We add this constraint to our tableau:

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 0 & 1 & \frac{10}{14} & -\frac{10}{14} & 0 & \frac{20}{14} \\ 1 & 0 & \frac{10}{14} & \frac{25}{14} & 0 & \frac{55}{14} \\ 0 & 0 & \frac{10}{14} & \frac{4}{14} & -1 & \frac{6}{14} \\ \mathbf{r}^\top & 0 & 0 & 5 & \frac{5}{2} & 0 & \frac{35}{2} \end{array}$$

Pivoting about the element (3, 3) gives

$$\begin{array}{cccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{b} \\ 0 & 1 & 0 & -1 & 1 & 1 \\ 1 & 0 & 0 & \frac{3}{2} & 1 & \frac{7}{2} \\ 0 & 0 & 1 & \frac{2}{5} & -\frac{7}{5} & \frac{3}{5} \\ \mathbf{r}^\top & 0 & 0 & 0 & \frac{1}{2} & 7 & \frac{29}{2} \end{array}$$

The corresponding optimal basic feasible solution is  $[7/2, 1, 3/5, 0, 0]^\top$ , which still does not satisfy the integer constraint.

Next, we construct the Gomory cut for the second row of the tableau:

$$\frac{1}{2}x_4 - x_6 = \frac{1}{2}$$

We add this constraint to our tableau to obtain

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
	0	1	0	-1	1	0	1
	1	0	0	$\frac{3}{2}$	1	0	$\frac{7}{2}$
	0	0	1	$\frac{2}{5}$	$-\frac{7}{5}$	0	$\frac{3}{5}$
	0	0	0	$\frac{1}{2}$	0	-1	$\frac{1}{2}$
$r^T$	0	0	0	$\frac{1}{2}$	7	0	$\frac{29}{2}$

Pivoting about (4, 4), we get

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
	0	1	0	0	1	-2	2
	1	0	0	0	1	3	2
	0	0	1	0	$-\frac{7}{5}$	$\frac{4}{5}$	$\frac{1}{5}$
	0	0	0	1	0	-2	1
$r^T$	0	0	0	0	7	1	14

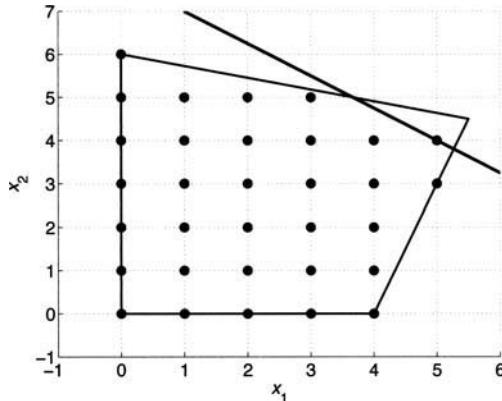
In this optimal basic feasible solution, the first two components are integers. Thus, we conclude that the solution to our ILP is  $[2, 2]^T$ , which agrees with the graphical solution in Figure 19.2. ■

In Example 19.5, the final solution to the LP problem after applying the Gomory cutting-plane method is not an integer vector. Only the first two components are integers, as these are the only two components in the original ILP problem. As pointed out earlier, the slack variables and variables introduced by the Gomory cuts are not constrained to be integers. However, if we are given an ILP problem with inequality constraints as in Example 19.5 but with only integer values in constraint data, then the slack variables and those introduced by the Gomory cuts are automatically integer valued (see also Exercise 19.9). We illustrate this in the following example.

**Example 19.6** Consider the following ILP problem:

$$\begin{aligned} & \text{maximize} && 3x_1 + 4x_2 \\ & \text{subject to} && 3x_1 - x_2 \leq 12 \\ & && 3x_1 + 11x_2 \leq 66 \\ & && x_1, x_2 \geq 0 \\ & && x_1, x_2 \in \mathbb{Z}. \end{aligned}$$

A graphical solution to this ILP problem is shown in Figure 19.3. As in Example 19.5, the solution is obtained by finding the straight line  $f = 3x_1 + 4x_2$  with largest  $f$  that passes through a feasible point with integer components. This point is  $[5, 4]^T$ .



**Figure 19.3** Graphical solution of the ILP problem in Example 19.6, where integer feasible solutions are marked with heavy dots.

We now solve the ILP problem above using the simplex method with Gomory cuts. We first represent the associated LP problem in standard form by introducing slack variables  $x_3$  and  $x_4$ . The initial tableau has the form

$$\begin{array}{ccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ 3 & -1 & 1 & 0 & 12 \\ 3 & 11 & 0 & 1 & 66 \\ \mathbf{c}^\top & -3 & -4 & 0 & 0 \end{array}$$

In this case there is an obvious initial basic feasible solution available, which allows us to initialize the simplex method to solve the problem. After two iterations of the simplex algorithm, the final tableau is

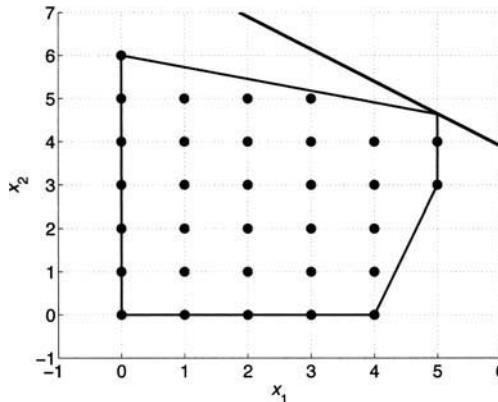
$$\begin{array}{ccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{b} \\ 1 & 0 & \frac{11}{36} & \frac{1}{36} & \frac{11}{2} \\ 0 & 1 & -\frac{1}{12} & \frac{1}{12} & \frac{9}{2} \\ \mathbf{r}^\top & 0 & 0 & \frac{7}{12} & \frac{5}{12} & \frac{69}{2} \end{array}$$

with optimal solution

$$\mathbf{x}^* = \left[ \frac{11}{2} \quad \frac{9}{2} \quad 0 \quad 0 \right]^\top.$$

Both basic components are noninteger. Let us construct a Gomory cut for the first basic component  $x_1^* = 11/2$ . From the first row of the tableau, the associated constraint equation is

$$x_1 + \frac{11}{36}x_3 + \frac{1}{36}x_4 = \frac{11}{2}.$$



**Figure 19.4** Graphical solution of the ILP in Example 19.6 after adding the constraint  $x_1 \leq 5$  to the original constraints.

If we apply the floor operator to this equation as explained before, we get an inequality constraint

$$x_1 \leq 5.$$

A graphical solution of the above problem after adding this inequality constraint to the original LP problem is shown in Figure 19.4. We can see that in this new problem, the first component of the optimal solution is an integer, but not the second. This means that a single Gomory cut will not suffice.

To continue with the Gomory procedure for the problem using the simplex method, we first write down the Gomory cut

$$\frac{11}{36}x_3 + \frac{1}{36}x_4 - x_5 = \frac{1}{2}.$$

We now obtain a new tableau by augmenting the previous tableau with the above constraint:

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$b$
1	0	$\frac{11}{36}$	$\frac{1}{36}$	0	$\frac{11}{2}$
0	1	$-\frac{1}{12}$	$\frac{1}{12}$	0	$\frac{9}{2}$
0	0	$\frac{11}{36}$	$\frac{1}{36}$	-1	$\frac{1}{2}$
$r^T$	0	0	$\frac{7}{12}$	$\frac{5}{12}$	0
					$\frac{69}{2}$

At this point, there is no obvious basic feasible solution. However, we can easily use the two-phase method. This yields

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$b$
1	0	0	0	1	5
0	1	0	$\frac{1}{11}$	$-\frac{3}{11}$	$\frac{51}{11}$
0	0	1	$\frac{1}{11}$	$-\frac{36}{11}$	$\frac{18}{11}$
$r^T$	0	0	0	$\frac{4}{11}$	$\frac{21}{11}$
					$\frac{369}{11}$

which has all nonnegative reduced cost coefficients. Hence, we obtain the optimal basic feasible solution

$$\mathbf{x}^* = \begin{bmatrix} 5 & \frac{51}{11} & \frac{18}{11} & 0 & 0 \end{bmatrix}^\top.$$

As expected, the second component does not satisfy the integer constraint.

Next, we write down the Gomory cut for the basic component  $x_2^* = 51/11$  using the numbers in the second row of the tableau:

$$\frac{1}{11}x_4 + \frac{8}{11}x_5 - x_6 = \frac{7}{11}.$$

Updating the tableau gives

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
1	0	0	0	1	0	5
0	1	0	$\frac{1}{11}$	$-\frac{3}{11}$	0	$\frac{51}{11}$
0	0	1	$\frac{1}{11}$	$-\frac{36}{11}$	0	$\frac{18}{11}$
0	0	0	$\frac{1}{11}$	$\frac{8}{11}$	-1	$\frac{7}{11}$
$r^\top$	0	0	0	$\frac{4}{11}$	$\frac{21}{11}$	0
						$\frac{369}{11}$

Again, there is no obvious basic feasible solution. Applying the two-phase method gives

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$b$
1	0	0	$-\frac{1}{8}$	0	$\frac{11}{8}$	$\frac{33}{8}$
0	1	0	$\frac{1}{8}$	0	$-\frac{3}{8}$	$\frac{39}{8}$
0	0	1	$\frac{1}{2}$	0	$-\frac{9}{2}$	$\frac{9}{2}$
0	0	0	$\frac{1}{8}$	1	$-\frac{11}{8}$	$\frac{7}{8}$
$r^\top$	0	0	0	$\frac{1}{8}$	0	$\frac{21}{8}$
						$\frac{255}{8}$

The corresponding optimal basic feasible solution still does not satisfy the integer constraints; neither the first nor the second components are integer.

Next, we introduce the Gomory cut using the numbers in the second row of the previous tableau to obtain

$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$	$b$
1	0	0	$-\frac{1}{8}$	0	$\frac{11}{8}$	0	$\frac{33}{8}$
0	1	0	$\frac{1}{8}$	0	$-\frac{3}{8}$	0	$\frac{39}{8}$
0	0	1	$\frac{1}{2}$	0	$-\frac{9}{2}$	0	$\frac{9}{2}$
0	0	0	$\frac{1}{8}$	1	$-\frac{11}{8}$	0	$\frac{7}{8}$
0	0	0	$\frac{1}{8}$	0	$\frac{5}{8}$	-1	$\frac{7}{8}$
$r^\top$	0	0	0	$\frac{1}{8}$	0	$\frac{21}{8}$	0
							$\frac{255}{8}$

Applying the two-phase method again gives

$$\begin{array}{ccccccccc} \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \mathbf{a}_4 & \mathbf{a}_5 & \mathbf{a}_6 & \mathbf{a}_7 & \mathbf{b} \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 5 \\ 0 & 1 & 0 & 0 & -\frac{1}{2} & 0 & \frac{1}{2} & 4 \\ 0 & 0 & 1 & 0 & -\frac{7}{2} & 0 & \frac{1}{2} & 1 \\ 0 & 0 & 0 & 1 & \frac{5}{2} & 0 & -\frac{11}{2} & 7 \\ 0 & 0 & 0 & 0 & -\frac{1}{2} & 1 & -\frac{1}{2} & 0 \\ \mathbf{r}^\top & 0 & 0 & 0 & 0 & 1 & 0 & 2 & 31 \end{array}$$

(Note that this basic feasible solution is degenerate—the corresponding basis is not unique.) The corresponding optimal basic feasible solution is

$$\begin{bmatrix} 5 & 4 & 1 & 7 & 0 & 0 & 0 \end{bmatrix}^\top,$$

which satisfies the integer constraints. From this, we see that the integer optimal solution to the original ILP problem is  $[5, 4]^\top$ , which agrees with our graphical solution in Figure 19.3.

In this example, we note that the final solution to LP problem after introducing slack variables and using the Gomory cutting-plane method is an integer vector. The reason for this, in contrast with Example 19.5, is that the original ILP inequality constraint data has only integers. ■

A linear programming problem in which not all of the components are required to be integers is called a *mixed integer linear programming (MILP)* problem. Gomory cuts are also relevant to solving MILP problems. In fact, Example 19.5 is an instance of an MILP problem, because the slack variables in the standard form of the problem are not constrained to be integers. Moreover, the cutting-plane idea also has been applied to nonsimplex methods and nonlinear programming algorithms.

For other methods for solving ILPs, see [119].

## EXERCISES

- 19.1** Show that if  $\mathbf{A}$  is totally unimodular, then so is any submatrix of it.
- 19.2** Show that if  $\mathbf{A}$  is totally unimodular, then so is  $\mathbf{A}^\top$ .
- 19.3** Show that  $\mathbf{A}$  is totally unimodular if and only if  $[\mathbf{A}, \mathbf{I}]$  is totally unimodular. This result is stronger than Proposition 19.1.
- 19.4** Consider the matrix  $\mathbf{B}$  in the proof of Proposition 19.1:

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_{m-k,k} & \mathbf{I}_{m-k} \\ \mathbf{B}_{k,k} & \mathbf{O} \end{bmatrix}.$$

Show that  $\det \mathbf{B} = \pm \det \mathbf{B}_{k,k}$ .

**19.5** Consider the constraint

$$\begin{aligned}\mathbf{A}\mathbf{x} &\leq \mathbf{b}, \\ \mathbf{x} &\in \mathbb{Z}^n,\end{aligned}$$

where  $\mathbf{A}$  and  $\mathbf{b}$  contain only integers. Suppose that we introduce the slack-variable vector  $\mathbf{z}$  to obtain the equivalent constraint

$$\begin{aligned}[\mathbf{A}, \mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} &= \mathbf{b} \\ \mathbf{x} &\in \mathbb{Z}^n \\ \mathbf{z} &\geq \mathbf{0}.\end{aligned}$$

Show that if  $\mathbf{z}$  satisfies this constraint (with some  $\mathbf{x}$ ), then  $\mathbf{z}$  is an integer vector.

**19.6** Write a MATLAB program to generate Figures 19.1 and 19.2.

**19.7** Consider the constraint in standard form  $\mathbf{Ax} = \mathbf{b}$ . Suppose that we augment this with a Gomory cut to obtain

$$\bar{\mathbf{A}} \begin{bmatrix} \mathbf{x} \\ x_{n+1} \end{bmatrix} = \bar{\mathbf{b}}.$$

Let  $x_{n+1}$  satisfy this constraint with an integer vector  $\mathbf{x}$ . Show that  $x_{n+1}$  is an integer.

**19.8** Consider the ILP problem in standard form

$$\begin{aligned}&\text{minimize } \mathbf{c}^\top \mathbf{x} \\ &\text{subject to } \mathbf{Ax} = \mathbf{b} \\ &\quad \mathbf{x} \geq \mathbf{0} \\ &\quad \mathbf{x} \in \mathbb{Z}^n.\end{aligned}$$

Show that if we use the Gomory cutting-plane method with the simplex algorithm, then the final optimal basic feasible solution, including the variables introduced by the Gomory method, is an integer vector. (Use Exercise 19.7.)

**19.9** Consider the ILP problem

$$\begin{aligned}&\text{minimize } \mathbf{c}^\top \mathbf{x} \\ &\text{subject to } \mathbf{Ax} \leq \mathbf{b} \\ &\quad \mathbf{x} \geq \mathbf{0} \\ &\quad \mathbf{x} \in \mathbb{Z}^n.\end{aligned}$$

Suppose that we introduce slack variables to convert the problem into standard form, and we use the Gomory cutting-plane method with the simplex algorithm to solve the resulting problem. Show that the final optimal basic feasible solution, including the slack variables and the variables introduced by the Gomory method, is an integer vector. (Use Exercises 19.5 and 19.8.)

- 19.10** Use a graphical method to find an integer solution to the dual of the ILP problem in Example 19.5.



## **PART IV**

---

# **NONLINEAR CONSTRAINED OPTIMIZATION**

---



# CHAPTER 20

---

## PROBLEMS WITH EQUALITY CONSTRAINTS

---

### 20.1 Introduction

In this part we discuss methods for solving a class of nonlinear constrained optimization problems that can be formulated as

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && h_i(\mathbf{x}) = 0, \quad i = 1, \dots, m \\ & && g_j(\mathbf{x}) \leq 0, \quad j = 1, \dots, p, \end{aligned}$$

where  $\mathbf{x} \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g_j : \mathbb{R}^n \rightarrow \mathbb{R}$ , and  $m \leq n$ . In vector notation, the problem above can be represented in the following *standard form*:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0} \\ & && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ . As usual, we adopt the following terminology.

**Definition 20.1** Any point satisfying the constraints is called a *feasible point*. The set of all feasible points,

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\},$$

is called a *feasible set*. ■

Optimization problems of the above form are not new to us. Indeed, linear programming problems of the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

which we studied in Part III, are of this type.

As we remarked in Part II, there is no loss of generality by considering only minimization problems. For if we are confronted with a maximization problem, it can easily be transformed into the minimization problem by observing that

$$\text{maximize } f(\mathbf{x}) = \text{minimize } -f(\mathbf{x}).$$

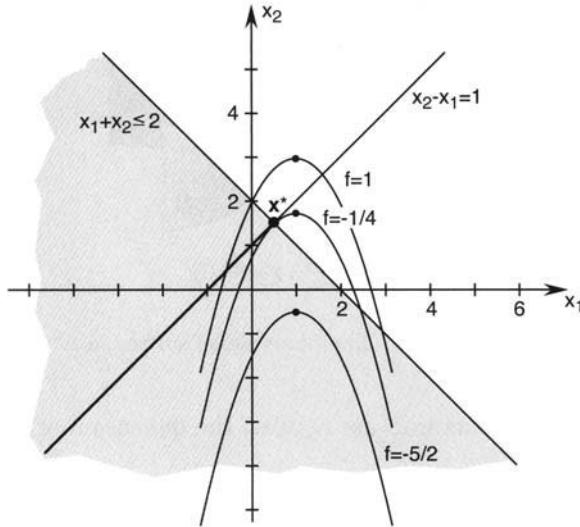
We illustrate the problems we study in this part by considering the following simple numerical example.

**Example 20.1** Consider the following optimization problem:

$$\begin{aligned} & \text{minimize} && (x_1 - 1)^2 + x_2 - 2 \\ & \text{subject to} && x_2 - x_1 = 1, \\ & && x_1 + x_2 \leq 2. \end{aligned}$$

This problem is already in the standard form given earlier, with  $f(x_1, x_2) = (x_1 - 1)^2 + x_2 - 2$ ,  $\mathbf{h}(x_1, x_2) = x_2 - x_1 - 1$ , and  $\mathbf{g}(x_1, x_2) = x_1 + x_2 - 2$ . This problem turns out to be simple enough to be solved graphically (see Figure 20.1). In the figure the set of points that satisfy the constraints (the feasible set) is marked by the heavy solid line. The inverted parabolae represent level sets of the objective function  $f$ —the lower the level set, the smaller the objective function value. Therefore, the solution can be obtained by finding the lowest-level set that intersects the feasible set. In this case, the minimizer lies on the level set with  $f = -1/4$ . The minimizer of the objective function is  $\mathbf{x}^* = [1/2, 3/2]^\top$ . ■

In the remainder of this chapter we discuss constrained optimization problems with only equality constraints. The general constrained optimization problem is discussed in the chapters to follow.



**Figure 20.1** Graphical solution to the problem in Example 20.1.

## 20.2 Problem Formulation

The class of optimization problems we analyze in this chapter is

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \end{aligned}$$

where  $\mathbf{x} \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{h} = [h_1, \dots, h_m]^\top$ , and  $m \leq n$ . We assume that the function  $\mathbf{h}$  is continuously differentiable, that is,  $\mathbf{h} \in \mathcal{C}^1$ .

We introduce the following definition.

**Definition 20.2** A point  $\mathbf{x}^*$  satisfying the constraints  $h_1(\mathbf{x}^*) = 0, \dots, h_m(\mathbf{x}^*) = 0$  is said to be a *regular point* of the constraints if the gradient vectors  $\nabla h_1(\mathbf{x}^*), \dots, \nabla h_m(\mathbf{x}^*)$  are linearly independent. ■

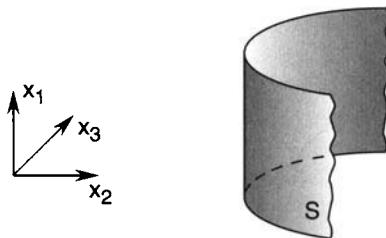
Let  $D\mathbf{h}(\mathbf{x}^*)$  be the Jacobian matrix of  $\mathbf{h} = [h_1, \dots, h_m]^\top$  at  $\mathbf{x}^*$ , given by

$$D\mathbf{h}(\mathbf{x}^*) = \begin{bmatrix} D h_1(\mathbf{x}^*) \\ \vdots \\ D h_m(\mathbf{x}^*) \end{bmatrix} = \begin{bmatrix} \nabla h_1(\mathbf{x}^*)^\top \\ \vdots \\ \nabla h_m(\mathbf{x}^*)^\top \end{bmatrix}.$$

Then,  $\mathbf{x}^*$  is regular if and only if  $\text{rank } D\mathbf{h}(\mathbf{x}^*) = m$  (i.e., the Jacobian matrix is of full rank).

The set of equality constraints  $h_1(\mathbf{x}) = 0, \dots, h_m(\mathbf{x}) = 0$ ,  $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ , describes a surface

$$S = \{\mathbf{x} \in \mathbb{R}^n : h_1(\mathbf{x}) = 0, \dots, h_m(\mathbf{x}) = 0\}.$$



$$S = \{[x_1, x_2, x_3]^\top : x_2 - x_3^2 = 0\}$$

**Figure 20.2** Two-dimensional surface in  $\mathbb{R}^3$ .

Assuming that the points in  $S$  are regular, the dimension of the surface  $S$  is  $n - m$ .

**Example 20.2** Let  $n = 3$  and  $m = 1$  (i.e., we are operating in  $\mathbb{R}^3$ ). Assuming that all points in  $S$  are regular, the set  $S$  is a two-dimensional surface. For example, let

$$h_1(\mathbf{x}) = x_2 - x_3^2 = 0.$$

Note that  $\nabla h_1(\mathbf{x}) = [0, 1, -2x_3]^\top$ , and hence for any  $\mathbf{x} \in \mathbb{R}^3$ ,  $\nabla h_1(\mathbf{x}) \neq \mathbf{0}$ . In this case,

$$\dim S = \dim \{\mathbf{x} : h_1(\mathbf{x}) = 0\} = n - m = 2.$$

See Figure 20.2 for a graphical illustration. ■

**Example 20.3** Let  $n = 3$  and  $m = 2$ . Assuming regularity, the feasible set  $S$  is a one-dimensional object (i.e., a curve in  $\mathbb{R}^3$ ). For example, let

$$\begin{aligned} h_1(\mathbf{x}) &= x_1, \\ h_2(\mathbf{x}) &= x_2 - x_3^2. \end{aligned}$$

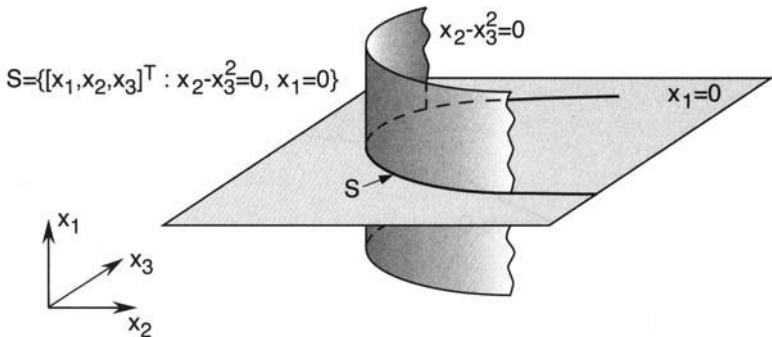
In this case,  $\nabla h_1(\mathbf{x}) = [1, 0, 0]^\top$  and  $\nabla h_2(\mathbf{x}) = [0, 1, -2x_3]^\top$ . Hence, the vectors  $\nabla h_1(\mathbf{x})$  and  $\nabla h_2(\mathbf{x})$  are linearly independent in  $\mathbb{R}^3$ . Thus,

$$\dim S = \dim \{\mathbf{x} : h_1(\mathbf{x}) = 0, h_2(\mathbf{x}) = 0\} = n - m = 1.$$

See Figure 20.3 for a graphical illustration. ■

### 20.3 Tangent and Normal Spaces

In this section we discuss the notion of a tangent space and normal space at a point on a surface. We begin by defining a *curve* on a surface  $S$ .



**Figure 20.3** One-dimensional surface in  $\mathbb{R}^3$ .

**Definition 20.3** A curve  $C$  on a surface  $S$  is a set of points  $\{\mathbf{x}(t) \in S : t \in (a, b)\}$ , continuously parameterized by  $t \in (a, b)$ ; that is,  $\mathbf{x} : (a, b) \rightarrow S$  is a continuous function. ■

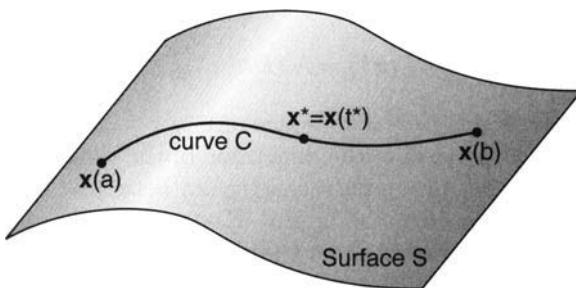
A graphical illustration of the definition of a curve is given in Figure 20.4. The definition of a curve implies that all the points on the curve satisfy the equation describing the surface. The curve  $C$  passes through a point  $\mathbf{x}^*$  if there exists  $t^* \in (a, b)$  such that  $\mathbf{x}(t^*) = \mathbf{x}^*$ .

Intuitively, we can think of a curve  $C = \{\mathbf{x}(t) : t \in (a, b)\}$  as the path traversed by a point  $\mathbf{x}$  traveling on the surface  $S$ . The position of the point at time  $t$  is given by  $\mathbf{x}(t)$ .

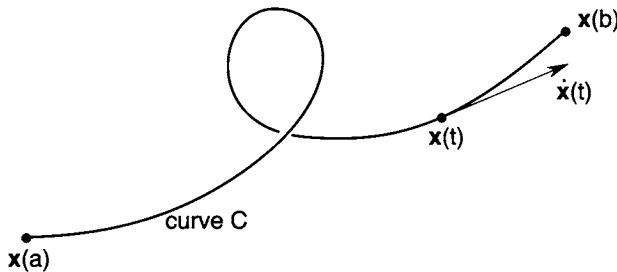
**Definition 20.4** The curve  $C = \{\mathbf{x}(t) : t \in (a, b)\}$  is differentiable if

$$\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}}{dt}(t) = \begin{bmatrix} \dot{x}_1(t) \\ \vdots \\ \dot{x}_n(t) \end{bmatrix}$$

exists for all  $t \in (a, b)$ .



**Figure 20.4** Curve on a surface.



**Figure 20.5** Geometric interpretation of the differentiability of a curve.

The curve  $C = \{\mathbf{x}(t) : t \in (a, b)\}$  is *twice differentiable* if

$$\ddot{\mathbf{x}}(t) = \frac{d^2\mathbf{x}}{dt^2}(t) = \begin{bmatrix} \ddot{x}_1(t) \\ \vdots \\ \ddot{x}_n(t) \end{bmatrix}$$

exists for all  $t \in (a, b)$ . ■

Note that both  $\dot{\mathbf{x}}(t)$  and  $\ddot{\mathbf{x}}(t)$  are  $n$ -dimensional vectors. We can think of  $\dot{\mathbf{x}}(t)$  and  $\ddot{\mathbf{x}}(t)$  as the velocity and acceleration, respectively, of a point traversing the curve  $C$  with position  $\mathbf{x}(t)$  at time  $t$ . The vector  $\dot{\mathbf{x}}(t)$  points in the direction of the instantaneous motion of  $\mathbf{x}(t)$ . Therefore, the vector  $\dot{\mathbf{x}}(t^*)$  is *tangent* to the curve  $C$  at  $\mathbf{x}^*$  (see Figure 20.5).

We are now ready to introduce the notions of a tangent space. For this recall the set

$$S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}\},$$

where  $\mathbf{h} \in \mathcal{C}^1$ . We think of  $S$  as a surface in  $\mathbb{R}^n$ .

**Definition 20.5** The *tangent space* at a point  $\mathbf{x}^*$  on the surface  $S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$  is the set  $T(\mathbf{x}^*) = \{\mathbf{y} : D\mathbf{h}(\mathbf{x}^*)\mathbf{y} = \mathbf{0}\}$ . ■

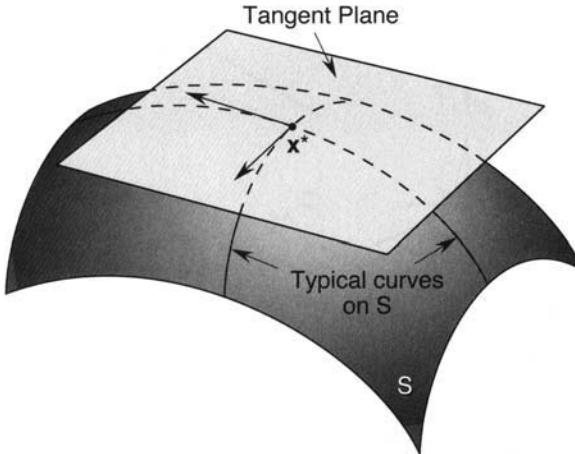
Note that the tangent space  $T(\mathbf{x}^*)$  is the nullspace of the matrix  $D\mathbf{h}(\mathbf{x}^*)$ :

$$T(\mathbf{x}^*) = \mathcal{N}(D\mathbf{h}(\mathbf{x}^*)).$$

The tangent space is therefore a subspace of  $\mathbb{R}^n$ .

Assuming that  $\mathbf{x}^*$  is regular, the dimension of the tangent space is  $n - m$ , where  $m$  is the number of equality constraints  $h_i(\mathbf{x}^*) = 0$ . Note that the tangent space passes through the origin. However, it is often convenient to picture the tangent space as a plane that passes through the point  $\mathbf{x}^*$ . For this, we define the *tangent plane* at  $\mathbf{x}^*$  to be the set

$$TP(\mathbf{x}^*) = T(\mathbf{x}^*) + \mathbf{x}^* = \{\mathbf{x} + \mathbf{x}^* : \mathbf{x} \in T(\mathbf{x}^*)\}.$$



**Figure 20.6** Tangent plane to the surface  $S$  at the point  $x^*$ .

Figure 20.6 illustrates the notion of a tangent plane, and Figure 20.7, the relationship between the tangent plane and the tangent space.

**Example 20.4** Let

$$S = \{\mathbf{x} \in \mathbb{R}^3 : h_1(\mathbf{x}) = x_1 = 0, h_2(\mathbf{x}) = x_1 - x_2 = 0\}.$$

Then,  $S$  is the  $x_3$ -axis in  $\mathbb{R}^3$  (see Figure 20.8). We have

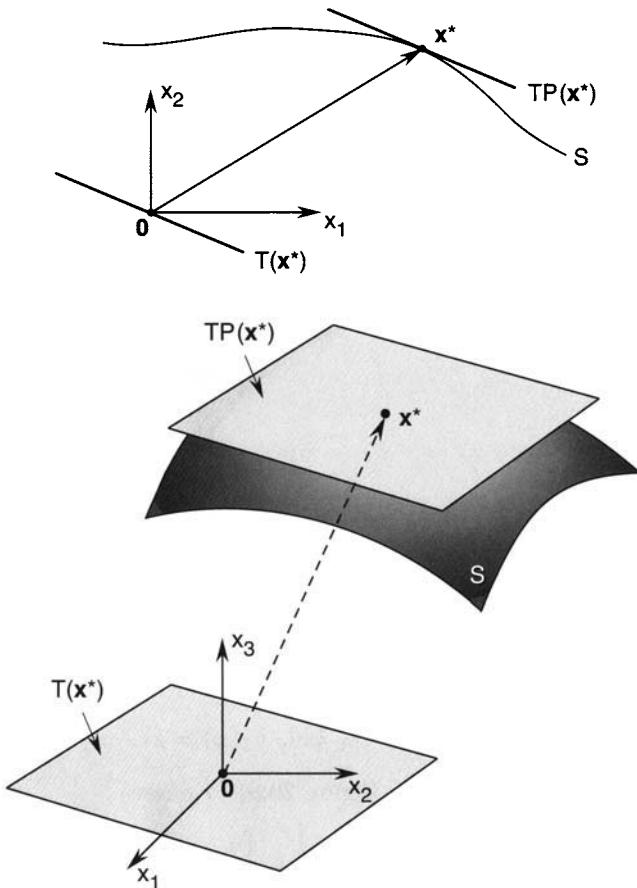
$$D\mathbf{h}(\mathbf{x}) = \begin{bmatrix} \nabla h_1(\mathbf{x})^\top \\ \nabla h_2(\mathbf{x})^\top \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & -1 & 0 \end{bmatrix}.$$

Because  $\nabla h_1$  and  $\nabla h_2$  are linearly independent when evaluated at any  $\mathbf{x} \in S$ , all the points of  $S$  are regular. The tangent space at an arbitrary point of  $S$  is

$$\begin{aligned} T(\mathbf{x}) &= \{\mathbf{y} : \nabla h_1(\mathbf{x})^\top \mathbf{y} = 0, \nabla h_2(\mathbf{x})^\top \mathbf{y} = 0\} \\ &= \left\{ \mathbf{y} : \begin{bmatrix} 1 & 0 & 0 \\ 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \mathbf{0} \right\} \\ &= \{[0, 0, \alpha]^\top : \alpha \in \mathbb{R}\} \\ &= \text{the } x_3\text{-axis in } \mathbb{R}^3. \end{aligned}$$

In this example, the tangent space  $T(\mathbf{x})$  at any point  $\mathbf{x} \in S$  is a one-dimensional subspace of  $\mathbb{R}^3$ . ■

Intuitively, we would expect the definition of the tangent space at a point on a surface to be the collection of all “tangent vectors” to the surface at that



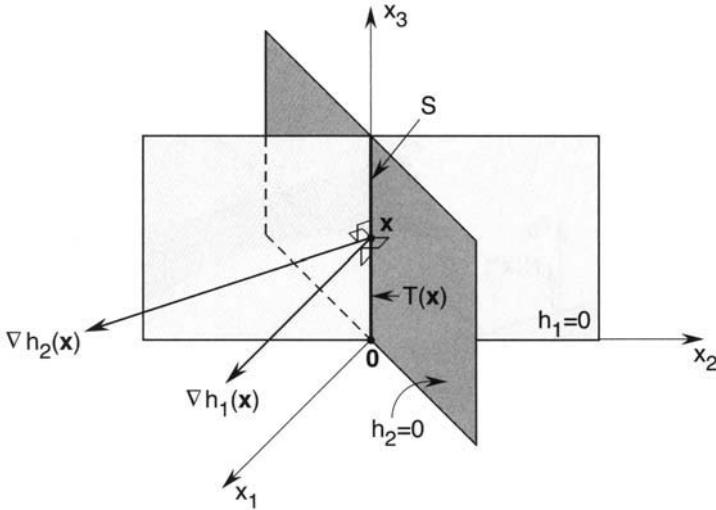
**Figure 20.7** Tangent spaces and planes in  $\mathbb{R}^2$  and  $\mathbb{R}^3$ .

point. We have seen that the derivative of a curve on a surface at a point is a tangent vector to the curve, and hence to the surface. The intuition above agrees with our definition whenever  $x^*$  is regular, as stated in the theorem below.

**Theorem 20.1** Suppose that  $x^* \in S$  is a regular point and  $T(x^*)$  is the tangent space at  $x^*$ . Then,  $y \in T(x^*)$  if and only if there exists a differentiable curve in  $S$  passing through  $x^*$  with derivative  $y$  at  $x^*$ .  $\square$

*Proof.*  $\Leftarrow$ : Suppose that there exists a curve  $\{x(t) : t \in (a, b)\}$  in  $S$  such that  $x(t^*) = x^*$  and  $\dot{x}(t^*) = y$  for some  $t^* \in (a, b)$ . Then,

$$h(x(t)) = 0$$



**Figure 20.8** The surface  $S = \{\mathbf{x} \in \mathbb{R}^3 : x_1 = 0, x_1 - x_2 = 0\}$ .

for all  $t \in (a, b)$ . If we differentiate the function  $\mathbf{h}(\mathbf{x}(t))$  with respect to  $t$  using the chain rule, we obtain

$$\frac{d}{dt} \mathbf{h}(\mathbf{x}(t)) = D\mathbf{h}(\mathbf{x}(t))\dot{\mathbf{x}}(t) = \mathbf{0}$$

for all  $t \in (a, b)$ . Therefore, at  $t^*$  we get

$$D\mathbf{h}(\mathbf{x}^*)\mathbf{y} = \mathbf{0},$$

and hence  $\mathbf{y} \in T(\mathbf{x}^*)$ .

⇒: To prove this, we need to use the implicit function theorem. We refer the reader to [88, p. 325]. ■

We now introduce the notion of a normal space.

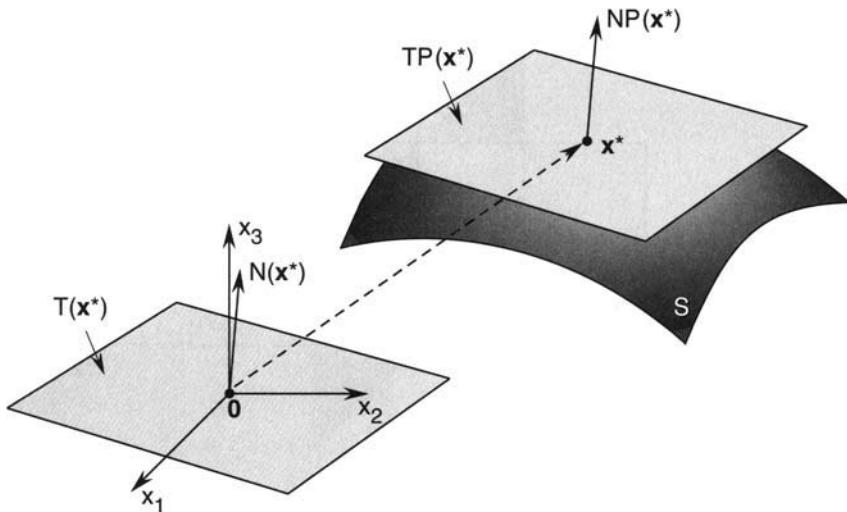
**Definition 20.6** The *normal space*  $N(\mathbf{x}^*)$  at a point  $\mathbf{x}^*$  on the surface  $S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$  is the set  $N(\mathbf{x}^*) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = D\mathbf{h}(\mathbf{x}^*)^\top \mathbf{z}, \mathbf{z} \in \mathbb{R}^m\}$ . ■

We can express the normal space  $N(\mathbf{x}^*)$  as

$$N(\mathbf{x}^*) = \mathcal{R}(D\mathbf{h}(\mathbf{x}^*)^\top),$$

that is, the range of the matrix  $D\mathbf{h}(\mathbf{x}^*)^\top$ . Note that the normal space  $N(\mathbf{x}^*)$  is the subspace of  $\mathbb{R}^n$  spanned by the vectors  $\nabla h_1(\mathbf{x}^*), \dots, \nabla h_m(\mathbf{x}^*)$ ; that is,

$$\begin{aligned} N(\mathbf{x}^*) &= \text{span}[\nabla h_1(\mathbf{x}^*), \dots, \nabla h_m(\mathbf{x}^*)] \\ &= \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = z_1 \nabla h_1(\mathbf{x}^*) + \dots + z_m \nabla h_m(\mathbf{x}^*), z_1, \dots, z_m \in \mathbb{R}\}. \end{aligned}$$



**Figure 20.9** Normal space in  $\mathbb{R}^3$ .

Note that the normal space contains the zero vector. Assuming that  $\mathbf{x}^*$  is regular, the dimension of the normal space  $N(\mathbf{x}^*)$  is  $m$ . As in the case of the tangent space, it is often convenient to picture the normal space  $N(\mathbf{x}^*)$  as passing through the point  $\mathbf{x}^*$  (rather than through the origin of  $\mathbb{R}^n$ ). For this, we define the *normal plane* at  $\mathbf{x}^*$  as the set

$$NP(\mathbf{x}^*) = N(\mathbf{x}^*) + \mathbf{x}^* = \{\mathbf{x} + \mathbf{x}^* \in \mathbb{R}^n : \mathbf{x} \in N(\mathbf{x}^*)\}.$$

Figure 20.9 illustrates the normal space and plane in  $\mathbb{R}^3$  (i.e.,  $n = 3$  and  $m = 1$ ).

We now show that the tangent space and normal space are orthogonal complements of each other (see Section 3.3).

**Lemma 20.1** *We have  $T(\mathbf{x}^*) = N(\mathbf{x}^*)^\perp$  and  $T(\mathbf{x}^*)^\perp = N(\mathbf{x}^*)$ .* □

*Proof.* By definition of  $T(\mathbf{x}^*)$ , we may write

$$T(\mathbf{x}^*) = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{x}^\top \mathbf{y} = 0 \text{ for all } \mathbf{x} \in N(\mathbf{x}^*)\}.$$

Hence, by definition of  $N(\mathbf{x}^*)$ , we have  $T(\mathbf{x}^*) = N(\mathbf{x}^*)^\perp$ . By Exercise 3.11 we also have  $T(\mathbf{x}^*)^\perp = N(\mathbf{x}^*)$ . ■

By Lemma 20.1, we can write  $\mathbb{R}^n$  as the direct sum decomposition (see Section 3.3):

$$\mathbb{R}^n = N(\mathbf{x}^*) \oplus T(\mathbf{x}^*);$$

that is, given any vector  $\mathbf{v} \in \mathbb{R}^n$ , there are unique vectors  $\mathbf{w} \in N(\mathbf{x}^*)$  and  $\mathbf{y} \in T(\mathbf{x}^*)$  such that

$$\mathbf{v} = \mathbf{w} + \mathbf{y}.$$

## 20.4 Lagrange Condition

In this section we present a first-order necessary condition for extremum problems with constraints. The result is the well-known *Lagrange's theorem*. To better understand the idea underlying this theorem, we first consider functions of two variables and only one equality constraint. Let  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$  be the constraint function. Recall that at each point  $\mathbf{x}$  of the domain, the gradient vector  $\nabla h(\mathbf{x})$  is orthogonal to the level set that passes through that point. Indeed, let us choose a point  $\mathbf{x}^* = [x_1^*, x_2^*]^\top$  such that  $h(\mathbf{x}^*) = 0$ , and assume that  $\nabla h(\mathbf{x}^*) \neq \mathbf{0}$ . The level set through the point  $\mathbf{x}^*$  is the set  $\{\mathbf{x} : h(\mathbf{x}) = 0\}$ . We then parameterize this level set in a neighborhood of  $\mathbf{x}^*$  by a curve  $\{\mathbf{x}(t)\}$ , that is, a continuously differentiable vector function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$  such that

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix}, \quad t \in (a, b), \quad \mathbf{x}^* = \mathbf{x}(t^*), \quad \dot{\mathbf{x}}(t^*) \neq \mathbf{0}, \quad t^* \in (a, b).$$

We can now show that  $\nabla h(\mathbf{x}^*)$  is orthogonal to  $\dot{\mathbf{x}}(t^*)$ . Indeed, because  $h$  is constant on the curve  $\{\mathbf{x}(t) : t \in (a, b)\}$ , we have that for all  $t \in (a, b)$ ,

$$h(\mathbf{x}(t)) = 0.$$

Hence, for all  $t \in (a, b)$ ,

$$\frac{d}{dt} h(\mathbf{x}(t)) = 0.$$

Applying the chain rule, we get

$$\frac{d}{dt} h(\mathbf{x}(t)) = \nabla h(\mathbf{x}(t))^\top \dot{\mathbf{x}}(t) = 0.$$

Therefore,  $\nabla h(\mathbf{x}^*)$  is orthogonal to  $\dot{\mathbf{x}}(t^*)$ .

Now suppose that  $\mathbf{x}^*$  is a minimizer of  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  on the set  $\{\mathbf{x} : h(\mathbf{x}) = 0\}$ . We claim that  $\nabla f(\mathbf{x}^*)$  is orthogonal to  $\dot{\mathbf{x}}(t^*)$ . To see this, it is enough to observe that the composite function of  $t$  given by

$$\phi(t) = f(\mathbf{x}(t))$$

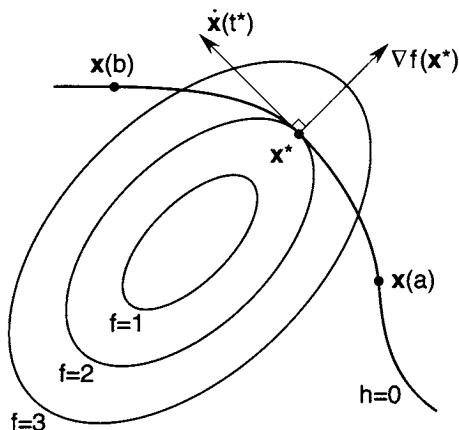
achieves a minimum at  $t^*$ . Consequently, the first-order necessary condition for the unconstrained extremum problem implies that

$$\frac{d\phi}{dt}(t^*) = 0.$$

Applying the chain rule yields

$$0 = \frac{d}{dt} \phi(t^*) = \nabla f(\mathbf{x}(t^*))^\top \dot{\mathbf{x}}(t^*) = \nabla f(\mathbf{x}^*)^\top \dot{\mathbf{x}}(t^*).$$

Thus,  $\nabla f(\mathbf{x}^*)$  is orthogonal to  $\dot{\mathbf{x}}(t^*)$ . The fact that  $\dot{\mathbf{x}}(t^*)$  is tangent to the curve  $\{\mathbf{x}(t)\}$  at  $\mathbf{x}^*$  means that  $\nabla f(\mathbf{x}^*)$  is orthogonal to the curve at  $\mathbf{x}^*$  (see Figure 20.10).



**Figure 20.10** The gradient  $\nabla f(\mathbf{x}^*)$  is orthogonal to the curve  $\{\mathbf{x}(t)\}$  at the point  $\mathbf{x}^*$  that is a minimizer of  $f$  on the curve.

Recall that  $\nabla h(\mathbf{x}^*)$  is also orthogonal to  $\dot{\mathbf{x}}(t^*)$ . Therefore, the vectors  $\nabla h(\mathbf{x}^*)$  and  $\nabla f(\mathbf{x}^*)$  are parallel; that is,  $\nabla f(\mathbf{x}^*)$  is a scalar multiple of  $\nabla h(\mathbf{x}^*)$ . The observations above allow us now to formulate *Lagrange's theorem* for functions of two variables with one constraint.

**Theorem 20.2 Lagrange's Theorem for  $n = 2, m = 1$ .** Let the point  $\mathbf{x}^*$  be a minimizer of  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  subject to the constraint  $h(\mathbf{x}) = 0$ ,  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ . Then,  $\nabla f(\mathbf{x}^*)$  and  $\nabla h(\mathbf{x}^*)$  are parallel. That is, if  $\nabla h(\mathbf{x}^*) \neq \mathbf{0}$ , then there exists a scalar  $\lambda^*$  such that

$$\nabla f(\mathbf{x}^*) + \lambda^* \nabla h(\mathbf{x}^*) = \mathbf{0}.$$

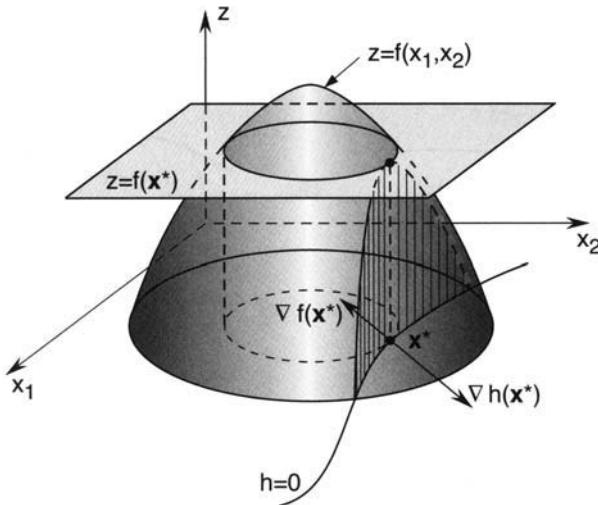
□

In Theorem 20.2, we refer to  $\lambda^*$  as the *Lagrange multiplier*. Note that the theorem also holds for maximizers. Figure 20.11 gives an illustration of Lagrange's theorem for the case where  $\mathbf{x}^*$  is a maximizer of  $f$  over the set  $\{\mathbf{x} : h(\mathbf{x}) = 0\}$ .

Lagrange's theorem provides a first-order necessary condition for a point to be a local minimizer. This condition, which we call the *Lagrange condition*, consists of two equations:

$$\begin{aligned}\nabla f(\mathbf{x}^*) + \lambda^* \nabla h(\mathbf{x}^*) &= \mathbf{0} \\ h(\mathbf{x}^*) &= 0.\end{aligned}$$

Note that the Lagrange condition is necessary but not sufficient. In Figure 20.12 we illustrate a variety of points where the Lagrange condition is



**Figure 20.11** Lagrange's theorem for  $n = 2, m = 1$ .

satisfied, including a case where the point is not an extremizer (neither a maximizer nor a minimizer).

We now generalize Lagrange's theorem for the case when  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ .

**Theorem 20.3 Lagrange's Theorem.** *Let  $\mathbf{x}^*$  be a local minimizer (or maximizer) of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ . Assume that  $\mathbf{x}^*$  is a regular point. Then, there exists  $\lambda^* \in \mathbb{R}^m$  such that*

$$Df(\mathbf{x}^*) + \lambda^{*\top} Dh(\mathbf{x}^*) = \mathbf{0}^\top.$$

□

*Proof.* We need to prove that

$$\nabla f(\mathbf{x}^*) = -D\mathbf{h}(\mathbf{x}^*)^\top \lambda^*$$

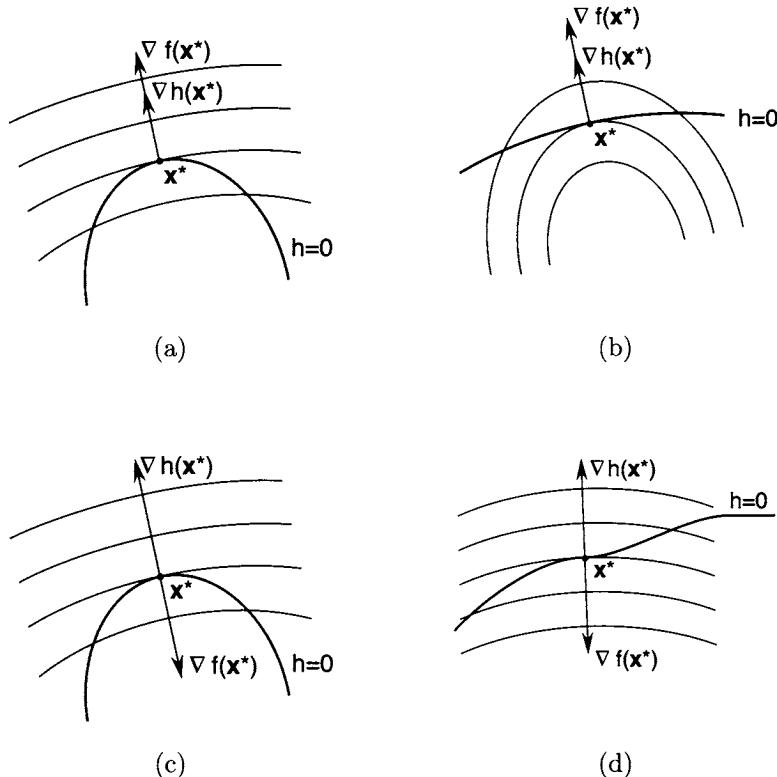
for some  $\lambda^* \in \mathbb{R}^m$ ; that is,  $\nabla f(\mathbf{x}^*) \in \mathcal{R}(D\mathbf{h}(\mathbf{x}^*)^\top) = N(\mathbf{x}^*)$ . But by Lemma 20.1,  $N(\mathbf{x}^*) = T(\mathbf{x}^*)^\perp$ . Therefore, it remains to show that  $\nabla f(\mathbf{x}^*) \in T(\mathbf{x}^*)^\perp$ .

We proceed as follows. Suppose that

$$\mathbf{y} \in T(\mathbf{x}^*).$$

Then, by Theorem 20.1, there exists a differentiable curve  $\{\mathbf{x}(t) : t \in (a, b)\}$  such that for all  $t \in (a, b)$ ,

$$\mathbf{h}(\mathbf{x}(t)) = \mathbf{0},$$



**Figure 20.12** Four examples where the Lagrange condition is satisfied: (a) maximizer, (b) minimizer, (c) minimizer, (d) not an extremizer. (Adapted from [120].)

and there exists  $t^* \in (a, b)$  satisfying

$$\mathbf{x}(t^*) = \mathbf{x}^*, \quad \dot{\mathbf{x}}(t^*) = \mathbf{y}.$$

Now consider the composite function  $\phi(t) = f(\mathbf{x}(t))$ . Note that  $t^*$  is a local minimizer of this function. By the first-order necessary condition for unconstrained local minimizers (see Theorem 6.1),

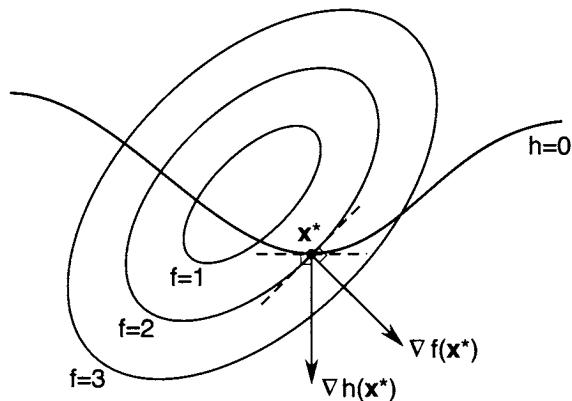
$$\frac{d\phi}{dt}(t^*) = 0.$$

Applying the chain rule yields

$$\frac{d\phi}{dt}(t^*) = Df(\mathbf{x}^*) \dot{\mathbf{x}}(t^*) = Df(\mathbf{x}^*) \mathbf{y} = \nabla f(\mathbf{x}^*)^\top \mathbf{y} = 0.$$

So all  $\mathbf{y} \in T(\mathbf{x}^*)$  satisfy

$$\nabla f(\mathbf{x}^*)^\top \mathbf{y} = 0;$$



**Figure 20.13** Example where the Lagrange condition does not hold.

that is,

$$\nabla f(\mathbf{x}^*) \in T(\mathbf{x}^*)^\perp.$$

This completes the proof. ■

Lagrange's theorem states that if  $\mathbf{x}^*$  is an extremizer, then the gradient of the objective function  $f$  can be expressed as a linear combination of the gradients of the constraints. We refer to the vector  $\lambda^*$  in Theorem 20.3 as the *Lagrange multiplier vector*, and its components as *Lagrange multipliers*.

From the proof of Lagrange's theorem, we see that a compact way to write the necessary condition is  $\nabla f(\mathbf{x}^*) \in N(\mathbf{x}^*)$ . If this condition fails, then  $\mathbf{x}^*$  cannot be an extremizer. This situation is illustrated in Figure 20.13.

Notice that regularity is stated as an assumption in Lagrange's theorem. This assumption plays an essential role, as illustrated in the following example.

**Example 20.5** Consider the following problem:

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && h(x) = 0, \end{aligned}$$

where  $f(x) = x$  and

$$h(x) = \begin{cases} x^2 & \text{if } x < 0 \\ 0 & \text{if } 0 \leq x \leq 1 \\ (x - 1)^2 & \text{if } x > 1. \end{cases}$$

The feasible set is evidently  $[0, 1]$ . Clearly,  $x^* = 0$  is a local minimizer. However,  $f'(x^*) = 1$  and  $h'(x^*) = 0$ . Therefore,  $x^*$  does not satisfy the

necessary condition in Lagrange's theorem. Note, however, that  $x^*$  is not a regular point, which is why Lagrange's theorem does not apply here. ■

It is convenient to introduce the *Lagrangian function*  $l : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ , given by

$$l(\mathbf{x}, \boldsymbol{\lambda}) \triangleq f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}).$$

The Lagrange condition for a local minimizer  $\mathbf{x}^*$  can be represented using the Lagrangian function as

$$Dl(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}^\top$$

for some  $\boldsymbol{\lambda}^*$ , where the derivative operation  $D$  is with respect to the entire argument  $[\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top$ . In other words, the necessary condition in Lagrange's theorem is equivalent to the first-order necessary condition for unconstrained optimization applied to the Lagrangian function.

To see the above, denote the derivative of  $l$  with respect to  $\mathbf{x}$  as  $D_x l$  and the derivative of  $l$  with respect to  $\boldsymbol{\lambda}$  as  $D_\lambda l$ . Then,

$$Dl(\mathbf{x}, \boldsymbol{\lambda}) = [D_x l(\mathbf{x}, \boldsymbol{\lambda}), D_\lambda l(\mathbf{x}, \boldsymbol{\lambda})].$$

Note that  $D_x l(\mathbf{x}, \boldsymbol{\lambda}) = Df(\mathbf{x}) + \boldsymbol{\lambda}^\top Dh(\mathbf{x})$  and  $D_\lambda l(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{h}(\mathbf{x})^\top$ . Therefore, Lagrange's theorem for a local minimizer  $\mathbf{x}^*$  can be stated as

$$\begin{aligned} D_x l(\mathbf{x}^*, \boldsymbol{\lambda}^*) &= \mathbf{0}^\top, \\ D_\lambda l(\mathbf{x}^*, \boldsymbol{\lambda}^*) &= \mathbf{0}^\top \end{aligned}$$

for some  $\boldsymbol{\lambda}^*$ , which is equivalent to

$$Dl(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}^\top.$$

In other words, the Lagrange condition can be expressed as  $Dl(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}^\top$ .

The Lagrange condition is used to find possible extremizers. This entails solving the equations

$$\begin{aligned} D_x l(\mathbf{x}, \boldsymbol{\lambda}) &= \mathbf{0}^\top, \\ D_\lambda l(\mathbf{x}, \boldsymbol{\lambda}) &= \mathbf{0}^\top. \end{aligned}$$

The above represents  $n + m$  equations in  $n + m$  unknowns. Keep in mind that the Lagrange condition is necessary but not sufficient; that is, a point  $\mathbf{x}^*$  satisfying the equations above need not be an extremizer.

**Example 20.6** Given a fixed area of cardboard, we wish to construct a closed cardboard box with maximum volume. We can formulate and solve this problem using the Lagrange condition. Denote the dimensions of the box with maximum volume by  $x_1$ ,  $x_2$ , and  $x_3$ , and let the given fixed area of cardboard be  $A$ . The problem can then be formulated as

$$\begin{aligned} &\text{maximize} && x_1 x_2 x_3 \\ &\text{subject to} && x_1 x_2 + x_2 x_3 + x_3 x_1 = \frac{A}{2}. \end{aligned}$$

We denote  $f(\mathbf{x}) = -x_1x_2x_3$  and  $h(\mathbf{x}) = x_1x_2 + x_2x_3 + x_3x_1 - A/2$ . We have  $\nabla f(\mathbf{x}) = -[x_2x_3, x_1x_3, x_1x_2]^\top$  and  $\nabla h(\mathbf{x}) = [x_2 + x_3, x_1 + x_3, x_1 + x_2]^\top$ . Note that all feasible points are regular in this case. By the Lagrange condition, the dimensions of the box with maximum volume satisfies

$$\begin{aligned} x_2x_3 - \lambda(x_2 + x_3) &= 0 \\ x_1x_3 - \lambda(x_1 + x_3) &= 0 \\ x_1x_2 - \lambda(x_1 + x_2) &= 0 \\ x_1x_2 + x_2x_3 + x_3x_1 &= \frac{A}{2}, \end{aligned}$$

where  $\lambda \in \mathbb{R}$ .

We now solve these equations. First, we show that that  $x_1$ ,  $x_2$ ,  $x_3$ , and  $\lambda$  are all nonzero. Suppose that  $x_1 = 0$ . By the constraints, we have  $x_2x_3 = A/2$ . However, the second and third equations in the Lagrange condition yield  $\lambda x_2 = \lambda x_3 = 0$ , which together with the first equation implies that  $x_2x_3 = 0$ . This contradicts the constraints. A similar argument applies to  $x_2$  and  $x_3$ .

Next, suppose that  $\lambda = 0$ . Then, the sum of the three Lagrange equations gives  $x_2x_3 + x_1x_3 + x_1x_2 = 0$ , which contradicts the constraints.

We now solve for  $x_1$ ,  $x_2$ , and  $x_3$  in the Lagrange equations. First, multiply the first equation by  $x_1$  and the second by  $x_2$ , and subtract one from the other. We arrive at  $x_3\lambda(x_1 - x_2) = 0$ . Because neither  $x_3$  nor  $\lambda$  can be zero (by part b), we conclude that  $x_1 = x_2$ . We similarly deduce that  $x_2 = x_3$ . From the constraint equation, we obtain  $x_1 = x_2 = x_3 = \sqrt{A/6}$ .

Notice that we have ignored the constraints that  $x_1$ ,  $x_2$ , and  $x_3$  are positive so that we can solve the problem using Lagrange's theorem. However, there is only one solution to the Lagrange equations, and the solution is positive. Therefore, if a solution exists for the problem with positivity constraints on the variables  $x_1$ ,  $x_2$ , and  $x_3$ , then this solution must necessarily be equal to the solution above obtained by ignoring the positivity constraints. ■

Next we provide an example with a quadratic objective function and a quadratic constraint.

**Example 20.7** Consider the problem of extremizing the objective function

$$f(\mathbf{x}) = x_1^2 + x_2^2$$

on the ellipse

$$\{[x_1, x_2]^\top : h(\mathbf{x}) = x_1^2 + 2x_2^2 - 1 = 0\}.$$

We have

$$\begin{aligned} \nabla f(\mathbf{x}) &= [2x_1, 2x_2]^\top, \\ \nabla h(\mathbf{x}) &= [2x_1, 4x_2]^\top. \end{aligned}$$

Thus,

$$D_{\mathbf{x}} l(\mathbf{x}, \lambda) = D_{\mathbf{x}} [f(\mathbf{x}) + \lambda h(\mathbf{x})] = [2x_1 + 2\lambda x_1, 2x_2 + 4\lambda x_2]$$

and

$$D_{\lambda} l(\mathbf{x}, \lambda) = h(\mathbf{x}) = x_1^2 + 2x_2^2 - 1.$$

Setting  $D_{\mathbf{x}} l(\mathbf{x}, \lambda) = \mathbf{0}^\top$  and  $D_{\lambda} l(\mathbf{x}, \lambda) = 0$ , we obtain three equations in three unknowns

$$\begin{aligned} 2x_1 + 2\lambda x_1 &= 0, \\ 2x_2 + 4\lambda x_2 &= 0, \\ x_1^2 + 2x_2^2 &= 1. \end{aligned}$$

All feasible points in this problem are regular. From the first of the equations above, we get either  $x_1 = 0$  or  $\lambda = -1$ . For the case where  $x_1 = 0$ , the second and third equations imply that  $\lambda = -1/2$  and  $x_2 = \pm 1/\sqrt{2}$ . For the case where  $\lambda = -1$ , the second and third equations imply that  $x_1 = \pm 1$  and  $x_2 = 0$ . Thus, the points that satisfy the Lagrange condition for extrema are

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0 \\ 1/\sqrt{2} \end{bmatrix}, \quad \mathbf{x}^{(2)} = \begin{bmatrix} 0 \\ -1/\sqrt{2} \end{bmatrix}, \quad \mathbf{x}^{(3)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{x}^{(4)} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}.$$

Because

$$f(\mathbf{x}^{(1)}) = f(\mathbf{x}^{(2)}) = \frac{1}{2}$$

and

$$f(\mathbf{x}^{(3)}) = f(\mathbf{x}^{(4)}) = 1$$

we conclude that if there are minimizers, then they are located at  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$ , and if there are maximizers, then they are located at  $\mathbf{x}^{(3)}$  and  $\mathbf{x}^{(4)}$ . It turns out that, indeed,  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  are minimizers and  $\mathbf{x}^{(3)}$  and  $\mathbf{x}^{(4)}$  are maximizers. This problem can be solved graphically, as illustrated in Figure 20.14. ■

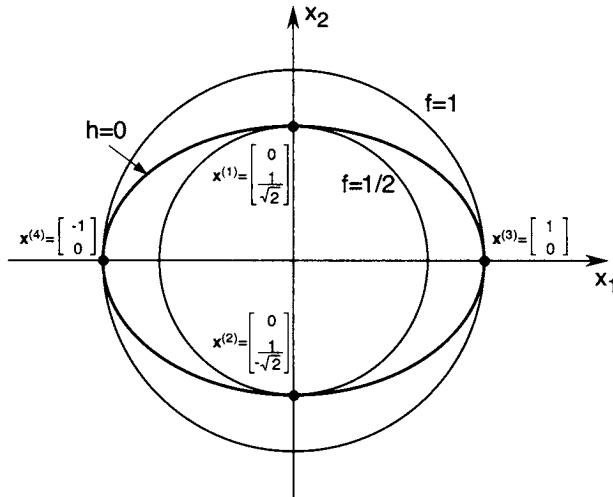
In the example above, both the objective function  $f$  and the constraint function  $h$  are quadratic functions. In the next example we take a closer look at a class of problems where both the objective function  $f$  and the constraint  $h$  are quadratic functions of  $n$  variables.

**Example 20.8** Consider the following problem:

$$\text{maximize } \frac{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}{\mathbf{x}^\top \mathbf{P} \mathbf{x}},$$

where  $\mathbf{Q} = \mathbf{Q}^\top \geq 0$  and  $\mathbf{P} = \mathbf{P}^\top > 0$ . Note that if a point  $\mathbf{x} = [x_1, \dots, x_n]^\top$  is a solution to the problem, then so is any nonzero scalar multiple of it,

$$t\mathbf{x} = [tx_1, \dots, tx_n]^\top, \quad t \neq 0.$$



**Figure 20.14** Graphical solution of the problem in Example 20.7.

Indeed,

$$\frac{(tx)^\top Q(tx)}{(tx)^\top P(tx)} = \frac{t^2 x^\top Q x}{t^2 x^\top P x} = \frac{x^\top Q x}{x^\top P x}.$$

Therefore, to avoid the multiplicity of solutions, we further impose the constraint

$$x^\top P x = 1.$$

The optimization problem becomes

$$\begin{aligned} & \text{maximize} && x^\top Q x \\ & \text{subject to} && x^\top P x = 1. \end{aligned}$$

Let us write

$$\begin{aligned} f(x) &= x^\top Q x, \\ h(x) &= 1 - x^\top P x. \end{aligned}$$

Any feasible point for this problem is regular (see Exercise 20.13). We now apply Lagrange's method. We first form the Lagrangian function

$$l(x, \lambda) = x^\top Q x + \lambda(1 - x^\top P x).$$

Applying the Lagrange condition yields

$$D_x l(x, \lambda) = 2x^\top Q - 2\lambda x^\top P = \mathbf{0}^\top,$$

$$D_\lambda l(x, \lambda) = 1 - x^\top P x = 0.$$

The first of the equations above can be represented as

$$\mathbf{Q}\mathbf{x} - \lambda\mathbf{P}\mathbf{x} = \mathbf{0}$$

or

$$(\lambda\mathbf{P} - \mathbf{Q})\mathbf{x} = \mathbf{0}.$$

This representation is possible because  $\mathbf{P} = \mathbf{P}^\top$  and  $\mathbf{Q} = \mathbf{Q}^\top$ . By assumption  $\mathbf{P} > 0$ , hence  $\mathbf{P}^{-1}$  exists. Premultiplying  $(\lambda\mathbf{P} - \mathbf{Q})\mathbf{x} = \mathbf{0}$  by  $\mathbf{P}^{-1}$ , we obtain

$$(\lambda\mathbf{I}_n - \mathbf{P}^{-1}\mathbf{Q})\mathbf{x} = \mathbf{0}$$

or, equivalently,

$$\mathbf{P}^{-1}\mathbf{Q}\mathbf{x} = \lambda\mathbf{x}.$$

Therefore, the solution, if it exists, is an eigenvector of  $\mathbf{P}^{-1}\mathbf{Q}$ , and the Lagrange multiplier is the corresponding eigenvalue. As usual, let  $\mathbf{x}^*$  and  $\lambda^*$  be the optimal solution. Because  $\mathbf{x}^{*\top}\mathbf{P}\mathbf{x}^* = 1$  and  $\mathbf{P}^{-1}\mathbf{Q}\mathbf{x}^* = \lambda^*\mathbf{x}^*$ , we have

$$\lambda^* = \mathbf{x}^{*\top}\mathbf{Q}\mathbf{x}^*.$$

Hence,  $\lambda^*$  is the maximum of the objective function, and therefore is, in fact, the maximal eigenvalue of  $\mathbf{P}^{-1}\mathbf{Q}$ . It is also called the maximal *generalized eigenvalue*. ■

In the problems above, we are able to find points that are candidates for extremizers of the given objective function subject to equality constraints. These critical points are the only candidates because they are the only points that satisfy the Lagrange condition. To classify such critical points as minimizers, maximizers, or neither, we need a stronger condition—possibly a necessary and sufficient condition. In the next section we discuss a second-order necessary condition and a second-order sufficient condition for minimizers.

## 20.5 Second-Order Conditions

We assume that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are twice continuously differentiable:  $f, \mathbf{h} \in \mathcal{C}^2$ . Let

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}) = f(\mathbf{x}) + \lambda_1 h_1(\mathbf{x}) + \cdots + \lambda_m h_m(\mathbf{x})$$

be the Lagrangian function. Let  $\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda})$  be the Hessian matrix of  $l(\mathbf{x}, \boldsymbol{\lambda})$  with respect to  $\mathbf{x}$ :

$$\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{F}(\mathbf{x}) + \lambda_1 \mathbf{H}_1(\mathbf{x}) + \cdots + \lambda_m \mathbf{H}_m(\mathbf{x}),$$

where  $\mathbf{F}(\mathbf{x})$  is the Hessian matrix of  $f$  at  $\mathbf{x}$  and  $\mathbf{H}_k(\mathbf{x})$  is the Hessian matrix of  $h_k$  at  $\mathbf{x}$ ,  $k = 1, \dots, m$ , given by

$$\mathbf{H}_k(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 h_k}{\partial x_1^2}(\mathbf{x}) & \cdots & \frac{\partial^2 h_k}{\partial x_n \partial x_1}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial^2 h_k}{\partial x_1 \partial x_n}(\mathbf{x}) & \cdots & \frac{\partial^2 h_k}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}.$$

We introduce the notation  $[\boldsymbol{\lambda} \mathbf{H}(\mathbf{x})]$ :

$$[\boldsymbol{\lambda} \mathbf{H}(\mathbf{x})] = \lambda_1 \mathbf{H}_1(\mathbf{x}) + \cdots + \lambda_m \mathbf{H}_m(\mathbf{x}).$$

Using the notation above, we can write

$$\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{F}(\mathbf{x}) + [\boldsymbol{\lambda} \mathbf{H}(\mathbf{x})].$$

**Theorem 20.4 Second-Order Necessary Conditions.** Let  $\mathbf{x}^*$  be a local minimizer of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ , and  $f, \mathbf{h} \in \mathcal{C}^2$ . Suppose that  $\mathbf{x}^*$  is regular. Then, there exists  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that:

1.  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) = \mathbf{0}^\top$ .
2. For all  $\mathbf{y} \in T(\mathbf{x}^*)$ , we have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} \geq 0$ . □

*Proof.* The existence of  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) = \mathbf{0}^\top$  follows from Lagrange's theorem. It remains to prove the second part of the result. Suppose that  $\mathbf{y} \in T(\mathbf{x}^*)$ ; that is,  $\mathbf{y}$  belongs to the tangent space to  $S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$  at  $\mathbf{x}^*$ . Because  $\mathbf{h} \in \mathcal{C}^2$ , following the argument of Theorem 20.1, there exists a twice-differentiable curve  $\{\mathbf{x}(t) : t \in (a, b)\}$  on  $S$  such that

$$\mathbf{x}(t^*) = \mathbf{x}^*, \quad \dot{\mathbf{x}}(t^*) = \mathbf{y}$$

for some  $t^* \in (a, b)$ . Observe that by assumption,  $t^*$  is a local minimizer of the function  $\phi(t) = f(\mathbf{x}(t))$ . From the second-order necessary condition for unconstrained minimization (see Theorem 6.2), we obtain

$$\frac{d^2\phi}{dt^2}(t^*) \geq 0.$$

Using the formula

$$\frac{d}{dt}(\mathbf{y}(t)^\top \mathbf{z}(t)) = \mathbf{z}(t)^\top \frac{d\mathbf{y}}{dt}(t) + \mathbf{y}(t)^\top \frac{d\mathbf{z}}{dt}(t)$$

and applying the chain rule yields

$$\begin{aligned} \frac{d^2\phi}{dt^2}(t^*) &= \frac{d}{dt}[Df(\mathbf{x}(t^*)) \dot{\mathbf{x}}(t^*)] \\ &= \dot{\mathbf{x}}(t^*)^\top \mathbf{F}(\mathbf{x}^*) \dot{\mathbf{x}}(t^*) + Df(\mathbf{x}^*) \ddot{\mathbf{x}}(t^*) \\ &= \mathbf{y}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{y} + Df(\mathbf{x}^*) \ddot{\mathbf{x}}(t^*) \geq 0. \end{aligned}$$

Because  $\mathbf{h}(\mathbf{x}(t)) = \mathbf{0}$  for all  $t \in (a, b)$ , we have

$$\frac{d^2}{dt^2} \boldsymbol{\lambda}^{*\top} \mathbf{h}(\mathbf{x}(t)) = 0.$$

Thus, for all  $t \in (a, b)$ ,

$$\begin{aligned}
\frac{d^2}{dt^2} \boldsymbol{\lambda}^{*\top} \mathbf{h}(\mathbf{x}(t)) &= \frac{d}{dt} \left[ \boldsymbol{\lambda}^{*\top} \frac{d}{dt} \mathbf{h}(\mathbf{x}(t)) \right] \\
&= \frac{d}{dt} \left[ \sum_{k=1}^m \lambda_k^* \frac{d}{dt} h_k(\mathbf{x}(t)) \right] \\
&= \frac{d}{dt} \left[ \sum_{k=1}^m \lambda_k^* D h_k(\mathbf{x}(t)) \dot{\mathbf{x}}(t) \right] \\
&= \sum_{k=1}^m \lambda_k^* \frac{d}{dt} (D h_k(\mathbf{x}(t)) \dot{\mathbf{x}}(t)) \\
&= \sum_{k=1}^m \lambda_k^* [\dot{\mathbf{x}}(t)^\top \mathbf{H}_k(\mathbf{x}(t)) \dot{\mathbf{x}}(t) + D h_k(\mathbf{x}(t)) \ddot{\mathbf{x}}(t)] \\
&= \dot{\mathbf{x}}^\top(t) [\boldsymbol{\lambda}^* \mathbf{H}(\mathbf{x}(t))] \dot{\mathbf{x}}(t) + \boldsymbol{\lambda}^{*\top} D \mathbf{h}(\mathbf{x}(t)) \ddot{\mathbf{x}}(t) \\
&= 0.
\end{aligned}$$

In particular, the above is true for  $t = t^*$ ; that is,

$$\mathbf{y}^\top [\boldsymbol{\lambda}^* \mathbf{H}(\mathbf{x}^*)] \mathbf{y} + \boldsymbol{\lambda}^{*\top} D \mathbf{h}(\mathbf{x}^*) \ddot{\mathbf{x}}(t^*) = 0.$$

Adding this equation to the inequality

$$\mathbf{y}^\top \mathbf{F}(\mathbf{x}^*) \mathbf{y} + D f(\mathbf{x}^*) \ddot{\mathbf{x}}(t^*) \geq 0$$

yields

$$\mathbf{y}^\top (\mathbf{F}(\mathbf{x}^*) + [\boldsymbol{\lambda}^* \mathbf{H}(\mathbf{x}^*)]) \mathbf{y} + (D f(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D \mathbf{h}(\mathbf{x}^*)) \ddot{\mathbf{x}}(t^*) \geq 0.$$

But, by Lagrange's theorem,  $D f(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D \mathbf{h}(\mathbf{x}^*) = \mathbf{0}^\top$ . Therefore,

$$\mathbf{y}^\top (\mathbf{F}(\mathbf{x}^*) + [\boldsymbol{\lambda}^* \mathbf{H}(\mathbf{x}^*)]) \mathbf{y} = \mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} \geq 0,$$

which proves the result. ■

Observe that  $\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda})$  plays a similar role as the Hessian matrix  $\mathbf{F}(\mathbf{x})$  of the objective function  $f$  did in the unconstrained minimization case. However, we now require that  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \geq 0$  only on  $T(\mathbf{x}^*)$  rather than on  $\mathbb{R}^n$ .

The conditions above are necessary, but not sufficient, for a point to be a local minimizer. We now present, without a proof, sufficient conditions for a point to be a strict local minimizer.

**Theorem 20.5 Second-Order Sufficient Conditions.** Suppose that  $f, \mathbf{h} \in C^2$  and there exists a point  $\mathbf{x}^* \in \mathbb{R}^n$  and  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that:

1.  $D f(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D \mathbf{h}(\mathbf{x}^*) = \mathbf{0}^\top$ .

2. For all  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ , we have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} > 0$ .

Then,  $\mathbf{x}^*$  is a strict local minimizer of  $f$  subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ .  $\square$

*Proof.* The interested reader can consult [88, p. 334] for a proof of this result.  $\blacksquare$

Theorem 20.5 states that if an  $\mathbf{x}^*$  satisfies the Lagrange condition, and  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  is positive definite on  $T(\mathbf{x}^*)$ , then  $\mathbf{x}^*$  is a strict local minimizer. A similar result to Theorem 20.5 holds for a strict local maximizer, the only difference being that  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  be negative definite on  $T(\mathbf{x}^*)$ . We illustrate this condition in the following example.

**Example 20.9** Consider the following problem:

$$\text{maximize } \frac{\mathbf{x}^\top \mathbf{Q}\mathbf{x}}{\mathbf{x}^\top \mathbf{P}\mathbf{x}},$$

where

$$\mathbf{Q} = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}.$$

As pointed out earlier, we can represent this problem in the equivalent form

$$\begin{aligned} &\text{maximize } \mathbf{x}^\top \mathbf{Q}\mathbf{x} \\ &\text{subject to } \mathbf{x}^\top \mathbf{P}\mathbf{x} = 1. \end{aligned}$$

The Lagrangian function for the transformed problem is given by

$$l(\mathbf{x}, \lambda) = \mathbf{x}^\top \mathbf{Q}\mathbf{x} + \lambda(1 - \mathbf{x}^\top \mathbf{P}\mathbf{x}).$$

The Lagrange condition yields

$$(\lambda\mathbf{I} - \mathbf{P}^{-1}\mathbf{Q})\mathbf{x} = \mathbf{0},$$

where

$$\mathbf{P}^{-1}\mathbf{Q} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}.$$

There are only two values of  $\lambda$  that satisfy  $(\lambda\mathbf{I} - \mathbf{P}^{-1}\mathbf{Q})\mathbf{x} = \mathbf{0}$ , namely, the eigenvalues of  $\mathbf{P}^{-1}\mathbf{Q}$ :  $\lambda_1 = 2$ ,  $\lambda_2 = 1$ . We recall from our previous discussion of this problem that the Lagrange multiplier corresponding to the solution is the maximum eigenvalue of  $\mathbf{P}^{-1}\mathbf{Q}$ , namely,  $\lambda^* = \lambda_1 = 2$ . The corresponding eigenvector is the maximizer—the solution to the problem. The eigenvector corresponding to the eigenvalue  $\lambda^* = 2$  satisfying the constraint  $\mathbf{x}^\top \mathbf{P}\mathbf{x} = 1$  is  $\pm\mathbf{x}^*$ , where

$$\mathbf{x}^* = \left[ \frac{1}{\sqrt{2}}, 0 \right]^\top.$$

At this point, all we have established is that the pairs  $(\pm \mathbf{x}^*, \lambda^*)$  satisfy the Lagrange condition. We now show that the points  $\pm \mathbf{x}^*$  are, in fact, strict local maximizers. We do this for the point  $\mathbf{x}^*$ . A similar procedure applies to  $-\mathbf{x}^*$ . We first compute the Hessian matrix of the Lagrangian function. We have

$$\mathbf{L}(\mathbf{x}^*, \lambda^*) = 2\mathbf{Q} - 2\lambda\mathbf{P} = \begin{bmatrix} 0 & 0 \\ 0 & -2 \end{bmatrix}.$$

The tangent space  $T(\mathbf{x}^*)$  to  $\{\mathbf{x} : 1 - \mathbf{x}^\top \mathbf{P}\mathbf{x} = 0\}$  is

$$\begin{aligned} T(\mathbf{x}^*) &= \{\mathbf{y} \in \mathbb{R}^2 : \mathbf{x}^{*\top} \mathbf{P}\mathbf{y} = 0\} \\ &= \{\mathbf{y} : [\sqrt{2}, 0]\mathbf{y} = 0\} \\ &= \{\mathbf{y} : \mathbf{y} = [0, a]^\top, a \in \mathbb{R}\}. \end{aligned}$$

Note that for each  $\mathbf{y} \in T(\mathbf{x}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ ,

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda^*) \mathbf{y} = [0, a] \begin{bmatrix} 0 & 0 \\ 0 & -2 \end{bmatrix} \begin{bmatrix} 0 \\ a \end{bmatrix} = -2a^2 < 0.$$

Hence,  $\mathbf{L}(\mathbf{x}^*, \lambda^*) < 0$  on  $T(\mathbf{x}^*)$ , and thus  $\mathbf{x}^* = [1/\sqrt{2}, 0]^\top$  is a strict local maximizer. The same is true for the point  $-\mathbf{x}^*$ . Note that

$$\frac{\mathbf{x}^{*\top} \mathbf{Q} \mathbf{x}^*}{\mathbf{x}^{*\top} \mathbf{P} \mathbf{x}^*} = 2,$$

which, as expected, is the value of the maximal eigenvalue of  $\mathbf{P}^{-1}\mathbf{Q}$ . Finally, we point out that any scalar multiple  $t\mathbf{x}^*$  of  $\mathbf{x}^*$ ,  $t \neq 0$ , is a solution to the original problem of maximizing  $\mathbf{x}^\top \mathbf{Q} \mathbf{x} / \mathbf{x}^\top \mathbf{P} \mathbf{x}$ . ■

## 20.6 Minimizing Quadratics Subject to Linear Constraints

Consider the problem

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \\ &\text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{Q} > 0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ . This problem is a special case of what is called a *quadratic programming problem* (the general form of a quadratic programming problem includes the constraint  $\mathbf{x} \geq \mathbf{0}$ ). Note that the constraint set contains an infinite number of points (see Section 2.3). We now show, using Lagrange's theorem, that there is a unique solution to the optimization problem above. Following that, we provide an example illustrating the application of this solution to an optimal control problem.

To solve the problem, we first form the Lagrangian function

$$l(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \boldsymbol{\lambda}^\top (\mathbf{b} - \mathbf{A} \mathbf{x}).$$

The Lagrange condition yields

$$D_{\mathbf{x}} l(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{x}^{*\top} \mathbf{Q} - \boldsymbol{\lambda}^{*\top} \mathbf{A} = \mathbf{0}^\top.$$

Rewriting, we get

$$\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{A}^\top \boldsymbol{\lambda}^*.$$

Premultiplying both sides of the above by  $\mathbf{A}$  gives

$$\mathbf{A} \mathbf{x}^* = \mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top \boldsymbol{\lambda}^*.$$

Using the fact that  $\mathbf{A} \mathbf{x}^* = \mathbf{b}$ , and noting that  $\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top$  is invertible because  $\mathbf{Q} > 0$  and  $\text{rank } \mathbf{A} = m$ , we can solve for  $\boldsymbol{\lambda}^*$  to obtain

$$\boldsymbol{\lambda}^* = (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} \mathbf{b}.$$

Therefore, we obtain

$$\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{A}^\top (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} \mathbf{b}.$$

The point  $\mathbf{x}^*$  is the only candidate for a minimizer. To establish that  $\mathbf{x}^*$  is indeed a minimizer, we verify that  $\mathbf{x}^*$  satisfies the second-order sufficient conditions. For this, we first find the Hessian matrix of the Lagrangian function at  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ . We have

$$\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{Q},$$

which is positive definite. Thus, the point  $\mathbf{x}^*$  is a strict local minimizer. We will see in Chapter 22 that  $\mathbf{x}^*$  is, in fact, a global minimizer.

The special case where  $\mathbf{Q} = \mathbf{I}_n$ , the  $n \times n$  identity matrix, reduces to the problem considered in Section 12.3. Specifically, the problem in Section 12.3 is to minimize the norm  $\|\mathbf{x}\|$  subject to  $\mathbf{A} \mathbf{x} = \mathbf{b}$ . The objective function here is  $f(\mathbf{x}) = \|\mathbf{x}\|$ , which is not differentiable at  $\mathbf{x} = \mathbf{0}$ . This precludes the use of Lagrange's theorem because the theorem requires differentiability of the objective function. We can overcome this difficulty by considering an equivalent optimization problem:

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|\mathbf{x}\|^2 \\ &\text{subject to} && \mathbf{A} \mathbf{x} = \mathbf{b}. \end{aligned}$$

The objective function  $\|\mathbf{x}\|^2/2$  has the same minimizer as the previous objective function  $\|\mathbf{x}\|$ . Indeed, if  $\mathbf{x}^*$  is such that for all  $\mathbf{x} \in \mathbb{R}^n$  satisfying  $\mathbf{A} \mathbf{x} = \mathbf{b}$ ,  $\|\mathbf{x}^*\| \leq \|\mathbf{x}\|$ , then  $\|\mathbf{x}^*\|^2/2 \leq \|\mathbf{x}\|^2/2$ . The same is true for the converse. Because the problem of minimizing  $\|\mathbf{x}\|^2/2$  subject to  $\mathbf{A} \mathbf{x} = \mathbf{b}$  is simply the

problem considered above with  $\mathbf{Q} = \mathbf{I}_n$ , we easily deduce the solution to be  $\mathbf{x}^* = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{b}$ , which agrees with the solution in Section 12.3.

**Example 20.10** Consider the discrete-time linear system model

$$x_k = ax_{k-1} + bu_k, \quad k \geq 1,$$

with initial condition  $x_0$  given. We can think of  $\{x_k\}$  as a discrete-time signal that is controlled by an external input signal  $\{u_k\}$ . In the control literature,  $x_k$  is called the *state* at time  $k$ . For a given  $x_0$ , our goal is to choose the control signal  $\{u_k\}$  so that the state remains “small” over a time interval  $[1, N]$ , but at the same time the control signal is “not too large.” To express the desire to keep the state  $\{x_k\}$  small, we choose the control sequence to minimize

$$\frac{1}{2} \sum_{i=1}^N x_i^2.$$

On the other hand, maintaining a control signal that is not too large, we minimize

$$\frac{1}{2} \sum_{i=1}^N u_i^2.$$

The two objectives above are conflicting in the sense that they cannot, in general, be achieved simultaneously—minimizing the first may result in a large control effort, while minimizing the second may result in large states. This is clearly a problem that requires compromise. One way to approach the problem is to minimize a weighted sum of the two functions above. Specifically, we can formulate the problem as

$$\text{minimize } \frac{1}{2} \sum_{i=1}^N (qx_i^2 + ru_i^2)$$

$$\text{subject to } x_k = ax_{k-1} + bu_k, \quad k = 1, \dots, N, \quad x_0 \text{ given,}$$

where the parameters  $q$  and  $r$  reflect the relative importance of keeping the state small versus keeping the control effort not too large. This problem is an instance of the *linear quadratic regulator (LQR) problem* (see, e.g., [15], [20], [85], [86], or [99]). Combining the two conflicting objectives of keeping the state small while keeping the control effort small is an instance of the *weighted sum* approach (see Section 24.4).

To solve the problem above, we can rewrite it as a quadratic programming problem. Define

$$\begin{aligned} \mathbf{Q} &= \begin{bmatrix} q\mathbf{I}_N & \mathbf{O} \\ \mathbf{O} & r\mathbf{I}_N \end{bmatrix}, \\ \mathbf{A} &= \begin{bmatrix} 1 & \dots & 0 & -b & \dots & 0 \\ -a & 1 & \vdots & -b & \vdots & \\ \ddots & \ddots & \vdots & & \ddots & \\ 0 & -a & 1 & 0 & \dots & -b \end{bmatrix}, \\ \mathbf{b} &= \begin{bmatrix} ax_0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{z} = [x_1, \dots, x_N, u_1, \dots, u_N]^\top. \end{aligned}$$

With these definitions, the problem reduces to the previously considered quadratic programming problem,

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \mathbf{z}^\top \mathbf{Q} \mathbf{z} \\ &\text{subject to} \quad \mathbf{A} \mathbf{z} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{Q}$  is  $2N \times 2N$ ,  $\mathbf{A}$  is  $N \times 2N$ , and  $\mathbf{b} \in \mathbb{R}^N$ . The solution is

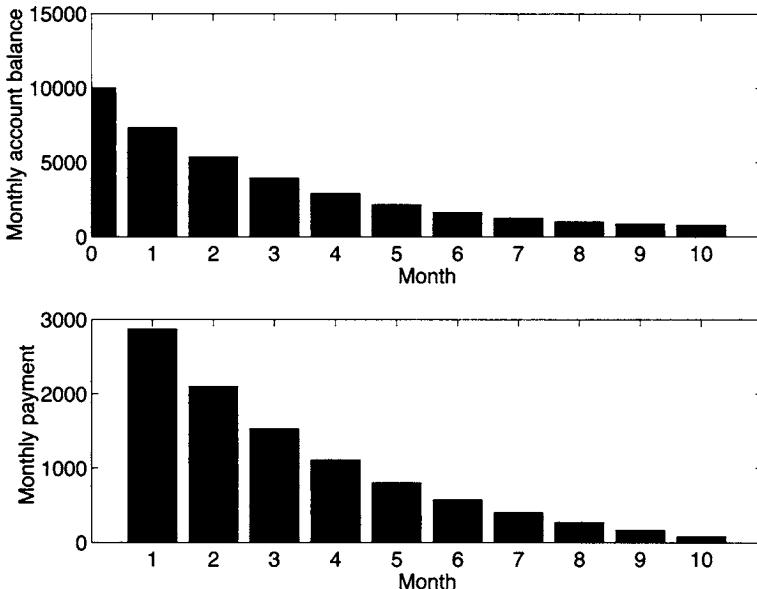
$$\mathbf{z}^* = \mathbf{Q}^{-1} \mathbf{A}^\top (\mathbf{A} \mathbf{Q}^{-1} \mathbf{A}^\top)^{-1} \mathbf{b}.$$

The first  $N$  components of  $\mathbf{z}^*$  represent the optimal state signal in the interval  $[1, N]$ , whereas the second  $N$  components represent the optimal control signal.

In practice, computation of the matrix inverses in the formula for  $\mathbf{z}^*$  above may be too costly. There are other ways to tackle the problem by exploiting its special structure. This is the study of *optimal control* (see, e.g., [15], [20], [85], [86], or [99]). ■

The following example illustrates an application of the above discussion.

**Example 20.11 Credit-Card Holder Dilemma.** Suppose that we currently have a credit-card debt of \$10,000. Credit-card debts are subject to a monthly interest rate of 2%, and the account balance is increased by the interest amount every month. Each month we have the option of reducing the account balance by contributing a payment to the account. Over the next 10 months, we plan to contribute a payment every month in such a way as to minimize the overall debt level while minimizing the hardship of making monthly payments.



**Figure 20.15** Plots for Example 20.11 with  $q = 1$  and  $r = 10$ .

We solve our problem using the LQR framework described in Example 20.10. Let the current time be 0,  $x_k$  the account balance at the end of month  $k$ , and  $u_k$  our payment in month  $k$ . We have

$$x_k = 1.02x_{k-1} - u_k, \quad k = 1, \dots, 10;$$

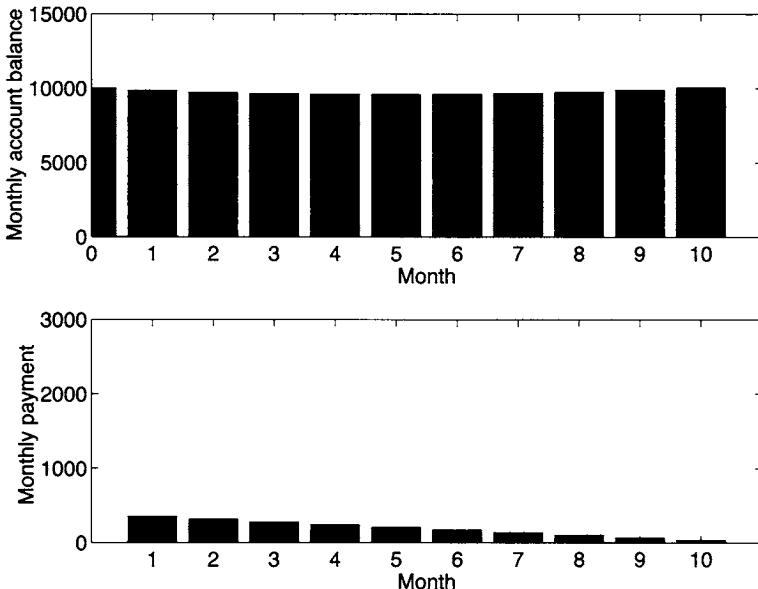
that is, the account balance in a given month is equal to the account balance in the previous month plus the monthly interest on that balance minus our payment that month. Our optimization problem is then

$$\text{minimize} \quad \frac{1}{2} \sum_{i=1}^{10} (qx_i^2 + ru_i^2)$$

$$\text{subject to} \quad x_k = 1.02x_{k-1} - u_k, \quad k = 1, \dots, 10, \quad x_0 = 10,000,$$

which is an instance of the LQR problem. The parameters  $q$  and  $r$  reflect our priority in trading off between debt reduction and hardship in making payments. The more anxious we are to reduce our debt, the larger the value of  $q$  relative to  $r$ . On the other hand, the more reluctant we are to make payments, the larger the value of  $r$  relative to  $q$ .

The solution to the problem above is given by the formula derived in Example 20.10. In Figure 20.15 we plot the monthly account balances and payments over the next 10 months using  $q = 1$  and  $r = 10$ . We can see here that our debt has been reduced to less than \$1000 after 10 months, but with a first



**Figure 20.16** Plots for Example 20.11 with  $q = 1$  and  $r = 300$ .

payment close to \$3000. If we feel that a payment of \$3000 is too high, then we can try to reduce this amount by increasing the value of  $r$  relative to  $q$ . However, going too far along these lines can lead to trouble. Indeed, if we use  $q = 1$  and  $r = 300$  (see Figure 20.16), although the monthly payments do not exceed \$400, the account balance is never reduced by much below \$10,000. In this case, the interest on the account balance eats up a significant portion of our monthly payments. In fact, our debt after 10 months will be higher than \$10,000. ■

For a treatment of optimization problems with quadratic objective functions, subject to linear or quadratic constraints, arising in communication and signal processing, see [105] and [106].

## EXERCISES

**20.1** Consider the following constraints on  $\mathbb{R}^2$ :

$$h(x_1, x_2) = (x_1 - 2)^2 = 0 \quad \text{and} \quad g(x_1, x_2) = (x_2 + 1)^3 \leq 0.$$

Find the set of feasible points. Are the feasible points regular? Justify your answer.

**20.2** Find local extremizers for the following optimization problems:

a. Minimize  $x_1^2 + 2x_1x_2 + 3x_2^2 + 4x_1 + 5x_2 + 6x_3$

subject to  $x_1 + 2x_2 = 3$

$4x_1 + 5x_3 = 6.$

b. Maximize  $4x_1 + x_2^2$

subject to  $x_1^2 + x_2^2 = 9.$

c. Maximize  $x_1x_2$

subject to  $x_1^2 + 4x_2^2 = 1.$

**20.3** Find minimizers and maximizers of the function

$$f(\mathbf{x}) = (\mathbf{a}^\top \mathbf{x})(\mathbf{b}^\top \mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^3,$$

subject to

$$x_1 + x_2 = 0$$

$$x_2 + x_3 = 0,$$

where

$$\mathbf{a} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}.$$

**20.4** Consider the problem

minimize  $f(\mathbf{x})$

subject to  $h(\mathbf{x}) = 0,$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ , and  $\nabla f(\mathbf{x}) = [x_1, x_1 + 4]^\top$ . Suppose that  $\mathbf{x}^*$  is an optimal solution and  $\nabla h(\mathbf{x}^*) = [1, 4]^\top$ . Find  $\nabla f(\mathbf{x}^*)$ .

**20.5** Consider the problem

minimize  $\|\mathbf{x} - \mathbf{x}_0\|^2$

subject to  $\|\mathbf{x}\|^2 = 9,$

where  $\mathbf{x}_0 = [1, \sqrt{3}]^\top$ .

- a. Find all points satisfying the Lagrange condition for the problem.

- b. Using second-order conditions, determine whether or not each of the points in part a is a local minimizer.

**20.6** We wish to construct a closed box with minimum surface area that encloses a volume of  $V$  cubic feet, where  $V > 0$ .

- Let  $a$ ,  $b$ , and  $c$  denote the dimensions of the box with minimum surface area (with volume  $V$ ). Derive the Lagrange condition that must be satisfied by  $a$ ,  $b$ , and  $c$ .
- What does it mean for a point  $\mathbf{x}^*$  to be a *regular* point in this problem? Is the point  $\mathbf{x}^* = [a, b, c]^\top$  a regular point?
- Find  $a$ ,  $b$ , and  $c$ .
- Does the point  $\mathbf{x}^* = [a, b, c]^\top$  found in part c satisfy the second-order sufficient condition?

**20.7** Find local extremizers of

- $f(x_1, x_2, x_3) = x_1^2 + 3x_2^2 + x_3$  subject to  $x_1^2 + x_2^2 + x_3^2 = 16$ .
- $f(x_1, x_2) = x_1^2 + x_2^2$  subject to  $3x_1^2 + 4x_1x_2 + 6x_2^2 = 140$ .

**20.8** Consider the problem

$$\begin{aligned} \text{minimize} \quad & 2x_1 + 3x_2 - 4, \quad x_1, x_2 \in \mathbb{R} \\ \text{subject to} \quad & x_1x_2 = 6. \end{aligned}$$

- Use Lagrange's theorem to find all possible local minimizers and maximizers.
- Use the second-order sufficient conditions to specify which points are strict local minimizers and which are strict local maximizers.
- Are the points in part b global minimizers or maximizers? Explain.

**20.9** Find all maximizers of the function

$$f(x_1, x_2) = \frac{18x_1^2 - 8x_1x_2 + 12x_2^2}{2x_1^2 + 2x_2^2}.$$

**20.10** Find all solutions to the problem

$$\begin{aligned} \text{maximize} \quad & \mathbf{x}^\top \begin{bmatrix} 3 & 4 \\ 0 & 3 \end{bmatrix} \mathbf{x} \\ \text{subject to} \quad & \|\mathbf{x}\|^2 = 1. \end{aligned}$$

**20.11** Consider a matrix  $\mathbf{A}$  with the property that  $\mathbf{A}^\top \mathbf{A}$  has eigenvalues ranging from 1 to 20 (i.e., the smallest eigenvalue is 1 and the largest is 20). Let  $\mathbf{x}$  be a vector such that  $\|\mathbf{x}\| = 1$ , and let  $\mathbf{y} = \mathbf{A}\mathbf{x}$ . Use Lagrange multiplier methods to find the range of values that  $\|\mathbf{y}\|$  can take.

*Hint:* What is the largest value that  $\|\mathbf{y}\|$  can take? What is the smallest value that  $\|\mathbf{y}\|$  can take?

**20.12** Consider a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Define the *induced 2-norm* of  $\mathbf{A}$ , denoted  $\|\mathbf{A}\|_2$ , to be the number

$$\|\mathbf{A}\|_2 = \max\{\|\mathbf{A}\mathbf{x}\| : \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\| = 1\},$$

where the norm  $\|\cdot\|$  on the right-hand side above is the usual Euclidean norm.

Suppose that the eigenvalues of  $\mathbf{A}^\top \mathbf{A}$  are  $\lambda_1, \dots, \lambda_n$  (ordered from largest to smallest). Use Lagrange's theorem to express  $\|\mathbf{A}\|_2$  in terms of the eigenvalues above (cf. Theorem 3.8).

**20.13** Let  $\mathbf{P} = \mathbf{P}^\top$  be a positive definite matrix. Show that any point  $\mathbf{x}$  satisfying  $1 - \mathbf{x}^\top \mathbf{P} \mathbf{x} = 0$  is a regular point.

**20.14** Consider the problem

$$\begin{aligned} & \text{maximize} && ax_1 + bx_2, \quad x_1, x_2 \in \mathbb{R} \\ & \text{subject to} && x_1^2 + x_2^2 = 2, \end{aligned}$$

where  $a, b \in \mathbb{R}$ . Show that if  $[1, 1]^\top$  is a solution to the problem, then  $a = b$ .

**20.15** Consider the problem

$$\begin{aligned} & \text{minimize} && x_1 x_2 - 2x_1, \quad x_1, x_2 \in \mathbb{R} \\ & \text{subject to} && x_1^2 - x_2^2 = 0. \end{aligned}$$

- a. Apply Lagrange's theorem directly to the problem to show that if a solution exists, it must be either  $[1, 1]^\top$  or  $[-1, 1]^\top$ .
- b. Use the second-order necessary conditions to show that  $[-1, 1]^\top$  cannot possibly be the solution.
- c. Use the second-order sufficient conditions to show that  $[1, 1]^\top$  is a strict local minimizer.

**20.16** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m \leq n$ ,  $\text{rank } \mathbf{A} = m$ , and  $\mathbf{x}_0 \in \mathbb{R}^n$ . Let  $\mathbf{x}^*$  be the point on the nullspace of  $\mathbf{A}$  that is closest to  $\mathbf{x}_0$  (in the sense of Euclidean norm).

- a. Show that  $\mathbf{x}^*$  is orthogonal to  $\mathbf{x}^* - \mathbf{x}_0$ .

- b. Find a formula for  $\mathbf{x}^*$  in terms of  $\mathbf{A}$  and  $\mathbf{x}_0$ .

**20.17** Consider the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \\ & \text{subject to} && \mathbf{Cx} = \mathbf{d}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m > n$ ,  $\mathbf{C} \in \mathbb{R}^{p \times n}$ ,  $p < n$ , and both  $\mathbf{A}$  and  $\mathbf{C}$  are of full rank. We wish to find an expression for the solution (in terms of  $\mathbf{A}$ ,  $\mathbf{b}$ ,  $\mathbf{C}$ , and  $\mathbf{d}$ ).

- a. Apply Lagrange's theorem to solve this problem.
- b. As an alternative, rewrite the given optimization problem in the form of a quadratic programming problem and apply the formula in Section 20.6 to obtain the solution.

**20.18** Consider the problem of minimizing a general quadratic function subject to a linear constraint:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Qx} - \mathbf{c}^\top \mathbf{x} + d \\ & \text{subject to} && \mathbf{Ax} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ , and  $d$  is a constant. Derive a closed-form solution to the problem.

**20.19** Let  $\mathbf{L}$  be an  $n \times n$  real symmetric matrix, and let  $\mathcal{M}$  be a subspace of  $\mathbb{R}^n$  with dimension  $m < n$ . Let  $\{\mathbf{b}_1, \dots, \mathbf{b}_m\} \subset \mathbb{R}^n$  be a basis for  $\mathcal{M}$ , and let  $\mathbf{B}$  be the  $n \times m$  matrix with  $\mathbf{b}_i$  as the  $i$ th column. Let  $\mathbf{L}_{\mathcal{M}}$  be the  $m \times m$  matrix defined by  $\mathbf{L}_{\mathcal{M}} = \mathbf{B}^\top \mathbf{LB}$ . Show that  $\mathbf{L}$  is positive semidefinite (definite) on  $\mathcal{M}$  if and only if  $\mathbf{L}_{\mathcal{M}}$  is positive semidefinite (definite).

*Note:* This result is useful for checking that the Hessian of the Lagrangian function at a point is positive definite on the tangent space at that point.

**20.20** Consider the sequence  $\{x_k\}$ ,  $x_k \in \mathbb{R}$ , generated by the recursion

$$x_{k+1} = ax_k + bu_k, \quad k \geq 0 \quad (a, b \in \mathbb{R}, a, b \neq 0),$$

where  $u_0, u_1, u_2, \dots$  is a sequence of “control inputs,” and the initial condition  $x_0 \neq 0$  is given. The recursion above is also called a *discrete-time linear system*. We wish to find values of control inputs  $u_0$  and  $u_1$  such that  $x_2 = 0$ , and the average input energy  $(u_0^2 + u_1^2)/2$  is minimized. Denote the optimal inputs by  $u_0^*$  and  $u_1^*$ .

- a. Find expressions for  $u_0^*$  and  $u_1^*$  in terms of  $a$ ,  $b$ , and  $x_0$ .
- b. Use the second-order sufficient conditions to show that the point  $\mathbf{u}^* = [u_0^*, u_1^*]^\top$  in part a is a strict local minimizer.

**20.21** Consider the discrete-time linear system  $x_k = 2x_{k-1} + u_k$ ,  $k \geq 1$ , with  $x_0 = 1$ . Find the values of the control inputs  $u_1$  and  $u_2$  to minimize

$$x_2^2 + \frac{1}{2}u_1^2 + \frac{1}{3}u_2^2.$$

**20.22** Consider the discrete-time linear system  $x_{k+1} = x_k + 2u_k$ ,  $0 \leq k \leq 2$ , with  $x_0 = 3$ . Use the Lagrange multiplier approach to calculate the optimal control sequence  $\{u_0, u_1, u_2\}$  that transfers the initial state  $x_0$  to  $x_3 = 9$  while minimizing

$$\frac{1}{2} \sum_{k=0}^2 u_k^2.$$

# CHAPTER 21

---

## PROBLEMS WITH INEQUALITY CONSTRAINTS

---

### 21.1 Karush-Kuhn-Tucker Condition

In Chapter 20 we analyzed constrained optimization problems involving only equality constraints. In this chapter we discuss extremum problems that also involve inequality constraints. The treatment in this chapter parallels that of Chapter 20. In particular, as we shall see, problems with inequality constraints can also be treated using Lagrange multipliers.

We consider the following problem:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}, \\ & && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ , and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ . For the general problem above, we adopt the following definitions.

**Definition 21.1** An inequality constraint  $g_j(\mathbf{x}) \leq 0$  is said to be *active* at  $\mathbf{x}^*$  if  $g_j(\mathbf{x}^*) = 0$ . It is *inactive* at  $\mathbf{x}^*$  if  $g_j(\mathbf{x}^*) < 0$ . ■

By convention, we consider an equality constraint  $h_i(\mathbf{x}) = 0$  to be always active.

**Definition 21.2** Let  $\mathbf{x}^*$  satisfy  $\mathbf{h}(\mathbf{x}^*) = \mathbf{0}$ ,  $\mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}$ , and let  $J(\mathbf{x}^*)$  be the index set of active inequality constraints:

$$J(\mathbf{x}^*) \triangleq \{j : g_j(\mathbf{x}^*) = 0\}.$$

Then, we say that  $\mathbf{x}^*$  is a *regular point* if the vectors

$$\nabla h_i(\mathbf{x}^*), \quad \nabla g_j(\mathbf{x}^*), \quad 1 \leq i \leq m, \quad j \in J(\mathbf{x}^*)$$

are linearly independent. ■

We now prove a first-order necessary condition for a point to be a local minimizer. We call this condition the *Karush-Kuhn-Tucker (KKT) condition*. In the literature, this condition is sometimes also called the Kuhn-Tucker condition.

**Theorem 21.1 Karush-Kuhn-Tucker (KKT) Theorem.** Let  $f, \mathbf{h}, \mathbf{g} \in C^1$ . Let  $\mathbf{x}^*$  be a regular point and a local minimizer for the problem of minimizing  $f$  subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ . Then, there exist  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$  such that:

1.  $\boldsymbol{\mu}^* \geq \mathbf{0}$ .
2.  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = 0$ .

□

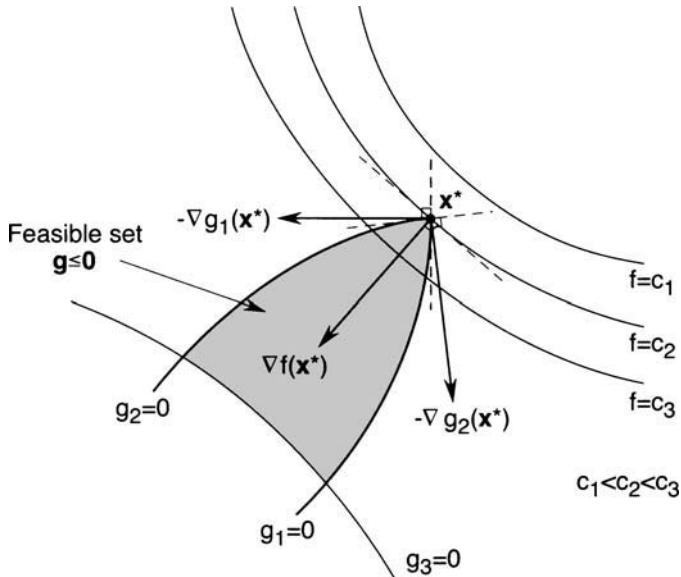
In Theorem 21.1, we refer to  $\boldsymbol{\lambda}^*$  as the Lagrange multiplier vector and  $\boldsymbol{\mu}^*$  as the Karush-Kuhn-Tucker (KKT) multiplier vector. We refer to their components as Lagrange multipliers and Karush-Kuhn-Tucker (KKT) multipliers, respectively.

Before proving this theorem, let us first discuss its meaning. Observe that  $\mu_j^* \geq 0$  (by condition 1) and  $g_j(\mathbf{x}^*) \leq 0$ . Therefore, the condition

$$\boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = \mu_1^* g_1(\mathbf{x}^*) + \cdots + \mu_p^* g_p(\mathbf{x}^*) = 0$$

implies that if  $g_j(\mathbf{x}^*) < 0$ , then  $\mu_j^* = 0$ ; that is, for all  $j \notin J(\mathbf{x}^*)$ , we have  $\mu_j^* = 0$ . In other words, the KKT multipliers  $\mu_j^*$  corresponding to inactive constraints are zero. The other KKT multipliers,  $\mu_i^*$ ,  $i \in J(\mathbf{x}^*)$ , are nonnegative; they may or may not be equal to zero.

**Example 21.1** A graphical illustration of the KKT theorem is given in Figure 21.1. In this two-dimensional example, we have only inequality constraints



**Figure 21.1** Illustration of the Karush-Kuhn-Tucker (KKT) theorem.

$g_j(\mathbf{x}) \leq 0, j = 1, 2, 3$ . Note that the point  $\mathbf{x}^*$  in the figure is indeed a minimizer. The constraint  $g_3(\mathbf{x}) \leq 0$  is inactive:  $g_3(\mathbf{x}^*) < 0$ ; hence  $\mu_3^* = 0$ . By the KKT theorem, we have

$$\nabla f(\mathbf{x}^*) + \mu_1^* \nabla g_1(\mathbf{x}^*) + \mu_2^* \nabla g_2(\mathbf{x}^*) = \mathbf{0},$$

or, equivalently,

$$\nabla f(\mathbf{x}^*) = -\mu_1^* \nabla g_1(\mathbf{x}^*) - \mu_2^* \nabla g_2(\mathbf{x}^*),$$

where  $\mu_1^* > 0$  and  $\mu_2^* > 0$ . It is easy to interpret the KKT condition graphically for this example. Specifically, we can see from Figure 21.1 that  $\nabla f(\mathbf{x}^*)$  must be a linear combination of the vectors  $-\nabla g_1(\mathbf{x}^*)$  and  $-\nabla g_2(\mathbf{x}^*)$  with positive coefficients. This is reflected exactly in the equation above, where the coefficients  $\mu_1^*$  and  $\mu_2^*$  are the KKT multipliers. ■

We apply the KKT condition in the same way that we apply any necessary condition. Specifically, we search for points satisfying the KKT condition and treat these points as candidate minimizers. To summarize, the KKT condition consists of five parts (three equations and two inequalities):

1.  $\boldsymbol{\mu}^* \geq \mathbf{0}$ .
2.  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\boldsymbol{\mu}^{*\top} g(\mathbf{x}^*) = 0$ .

$$4. \quad \mathbf{h}(\mathbf{x}^*) = \mathbf{0}.$$

$$5. \quad \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}.$$

We now prove the KKT theorem.

*Proof of the Karush-Kuhn-Tucker Theorem.* Let  $\mathbf{x}^*$  be a regular local minimizer of  $f$  on the set  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ . Then,  $\mathbf{x}^*$  is also a regular local minimizer of  $f$  on the set  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, g_j(\mathbf{x}) = 0, j \in J(\mathbf{x}^*)\}$  (see Exercise 21.16). Note that the latter constraint set involves only equality constraints. Therefore, from Lagrange's theorem, it follows that there exist vectors  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$  such that

$$Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top,$$

where for all  $j \notin J(\mathbf{x}^*)$ , we have  $\mu_j^* = 0$ . To complete the proof it remains to show that for all  $j \in J(\mathbf{x}^*)$ , we have  $\mu_j^* \geq 0$  (and hence for all  $j = 1, \dots, p$ , we have  $\mu_j^* \geq 0$ , i.e.,  $\boldsymbol{\mu}^* \geq \mathbf{0}$ ). We use a proof by contradiction. So suppose that there exists  $j \in J(\mathbf{x}^*)$  such that  $\mu_j^* < 0$ . Let  $\hat{S}$  and  $\hat{T}(\mathbf{x}^*)$  be the surface and tangent space defined by all other active constraints at  $\mathbf{x}^*$ . Specifically,

$$\hat{S} = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, g_i(\mathbf{x}) = 0, i \in J(\mathbf{x}^*), i \neq j\}$$

and

$$\hat{T}(\mathbf{x}^*) = \{\mathbf{y} : D\mathbf{h}(\mathbf{x}^*)\mathbf{y} = \mathbf{0}, Dg_i(\mathbf{x}^*)\mathbf{y} = 0, i \in J(\mathbf{x}^*), i \neq j\}.$$

We claim that by the regularity of  $\mathbf{x}^*$ , there exists  $\mathbf{y} \in \hat{T}(\mathbf{x}^*)$  such that

$$Dg_j(\mathbf{x}^*)\mathbf{y} \neq 0.$$

To see this, suppose that for all  $\mathbf{y} \in \hat{T}(\mathbf{x}^*)$ ,  $\nabla g_j(\mathbf{x}^*)^\top \mathbf{y} = Dg_j(\mathbf{x}^*)\mathbf{y} = 0$ . This implies that  $\nabla g_j(\mathbf{x}^*) \in \hat{T}(\mathbf{x}^*)^\perp$ . By Lemma 20.1, this, in turn, implies that

$$\nabla g_j(\mathbf{x}^*) \in \text{span}[\nabla h_k(\mathbf{x}^*), k = 1, \dots, m, \nabla g_i(\mathbf{x}^*), i \in J(\mathbf{x}^*), i \neq j].$$

But this contradicts the fact that  $\mathbf{x}^*$  is a regular point, which proves our claim. Without loss of generality, we assume that we have  $\mathbf{y}$  such that  $Dg_j(\mathbf{x}^*)\mathbf{y} < 0$ .

Consider the Lagrange condition, rewritten as

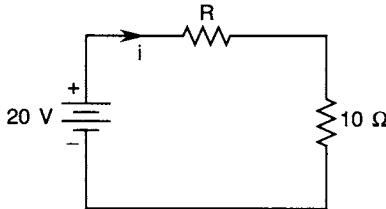
$$Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \mu_j^* Dg_j(\mathbf{x}^*) + \sum_{i \neq j} \mu_i^* Dg_i(\mathbf{x}^*) = \mathbf{0}^\top.$$

If we postmultiply the above by  $\mathbf{y}$  and use the fact that  $\mathbf{y} \in \hat{T}(\mathbf{x}^*)$ , we get

$$Df(\mathbf{x}^*)\mathbf{y} = -\mu_j^* Dg_j(\mathbf{x}^*)\mathbf{y}.$$

Because  $Dg_j(\mathbf{x}^*)\mathbf{y} < 0$  and we have assumed that  $\mu_j^* < 0$ , we have

$$Df(\mathbf{x}^*)\mathbf{y} < 0.$$



**Figure 21.2** Circuit in Example 21.2.

Because  $\mathbf{y} \in \hat{T}(\mathbf{x}^*)$ , by Theorem 20.1 we can find a differentiable curve  $\{\mathbf{x}(t) : t \in (a, b)\}$  on  $\hat{S}$  such that there exists  $t^* \in (a, b)$  with  $\mathbf{x}(t^*) = \mathbf{x}^*$  and  $\dot{\mathbf{x}}(t^*) = \mathbf{y}$ . Now,

$$\frac{d}{dt} f(\mathbf{x}(t^*)) = Df(\mathbf{x}^*)\mathbf{y} < 0.$$

The above means that there is a  $\delta > 0$  such that for all  $t \in (t^*, t^* + \delta]$ , we have

$$f(\mathbf{x}(t)) < f(\mathbf{x}(t^*)) = f(\mathbf{x}^*).$$

On the other hand,

$$\frac{d}{dt} g_j(\mathbf{x}(t^*)) = Dg_j(\mathbf{x}^*)\mathbf{y} < 0,$$

and for some  $\varepsilon > 0$  and all  $t \in [t^*, t^* + \varepsilon]$ , we have that  $g_j(\mathbf{x}(t)) \leq 0$ . Therefore, for all  $t \in (t^*, t^* + \min\{\delta, \varepsilon\}]$ , we have that  $g_j(\mathbf{x}(t)) \leq 0$  and  $f(\mathbf{x}(t)) < f(\mathbf{x}^*)$ . Because the points  $\mathbf{x}(t)$ ,  $t \in (t^*, t^* + \min\{\delta, \varepsilon\}]$ , are in  $\hat{S}$ , they are feasible points with lower objective function values than  $\mathbf{x}^*$ . This contradicts the assumption that  $\mathbf{x}^*$  is a local minimizer, which completes the proof. ■

**Example 21.2** Consider the circuit in Figure 21.2. Formulate and solve the KKT condition for the following problems.

- Find the value of the resistor  $R \geq 0$  such that the power absorbed by this resistor is maximized.
- Find the value of the resistor  $R \geq 0$  such that the power delivered to the  $10\text{-}\Omega$  resistor is maximized.

*Solution:*

- The power absorbed by the resistor  $R$  is  $p = i^2 R$ , where  $i = \frac{20}{10+R}$ . The optimization problem can be represented as

$$\begin{aligned} \text{minimize} \quad & -\frac{400R}{(10+R)^2} \\ \text{subject to} \quad & -R \leq 0. \end{aligned}$$

The derivative of the objective function is

$$-\frac{400(10+R)^2 - 800R(10+R)}{(10+R)^4} = -\frac{400(10-R)}{(10+R)^3}.$$

Thus, the KKT condition is

$$\begin{aligned} -\frac{400(10-R)}{(10+R)^3} - \mu &= 0, \\ \mu &\geq 0, \\ \mu R &= 0, \\ -R &\leq 0. \end{aligned}$$

We consider two cases. In the first case, suppose that  $\mu > 0$ . Then,  $R = 0$ . But this contradicts the first condition above. Now suppose that  $\mu = 0$ . Then, by the first condition, we have  $R = 10$ . Therefore, the only solution to the KKT condition is  $R = 10$ ,  $\mu = 0$ .

- b. The power absorbed by the  $10\text{-}\Omega$  resistor is  $p = i^2 10$ , where  $i = 20/(10+R)$ . The optimization problem can be represented as

$$\begin{aligned} \text{minimize} \quad & -\frac{4000}{(10+R)^2} \\ \text{subject to} \quad & -R \leq 0. \end{aligned}$$

The derivative of the objective function is

$$\frac{8000}{(10+R)^3}.$$

Thus, the KKT condition is

$$\begin{aligned} \frac{8000}{(10+R)^3} - \mu &= 0, \\ \mu &\geq 0, \\ \mu R &= 0, \\ -R &\leq 0. \end{aligned}$$

As before, we consider two cases. In the first case, suppose that  $\mu > 0$ . Then,  $R = 0$ , which is feasible. For the second case, suppose that  $\mu = 0$ . But this contradicts the first condition. Therefore, the only solution to the KKT condition is  $R = 0$ ,  $\mu = 8$ . ■

In the case when the objective function is to be maximized, that is, when the optimization problem has the form

$$\begin{aligned} \text{maximize} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & \mathbf{h}(\mathbf{x}) = \mathbf{0} \\ & \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

the KKT condition can be written as

1.  $\mu^* \geq \mathbf{0}$ .
2.  $-Df(\mathbf{x}^*) + \lambda^{*\top} Dh(\mathbf{x}^*) + \mu^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\mu^{*\top} g(\mathbf{x}^*) = 0$ .
4.  $h(\mathbf{x}^*) = \mathbf{0}$ .
5.  $g(\mathbf{x}^*) \leq \mathbf{0}$ .

The above is easily derived by converting the maximization problem above into a minimization problem, by multiplying the objective function by  $-1$ . It can be further rewritten as

1.  $\mu^* \leq \mathbf{0}$ .
2.  $Df(\mathbf{x}^*) + \lambda^{*\top} Dh(\mathbf{x}^*) + \mu^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\mu^{*\top} g(\mathbf{x}^*) = 0$ .
4.  $h(\mathbf{x}^*) = \mathbf{0}$ .
5.  $g(\mathbf{x}^*) \leq \mathbf{0}$ .

The form shown above is obtained from the preceding one by changing the signs of  $\mu^*$  and  $\lambda^*$  and multiplying condition 2 by  $-1$ .

We can similarly derive the KKT condition for the case when the inequality constraint is of the form  $g(\mathbf{x}) \geq \mathbf{0}$ . Specifically, consider the problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && h(\mathbf{x}) = \mathbf{0} \\ & && g(\mathbf{x}) \geq \mathbf{0}. \end{aligned}$$

We multiply the inequality constraint function by  $-1$  to obtain  $-g(\mathbf{x}) \leq \mathbf{0}$ . Thus, the KKT condition for this case is

1.  $\mu^* \geq \mathbf{0}$ .
2.  $Df(\mathbf{x}^*) + \lambda^{*\top} Dh(\mathbf{x}^*) - \mu^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\mu^{*\top} g(\mathbf{x}^*) = 0$ .
4.  $h(\mathbf{x}^*) = \mathbf{0}$ .
5.  $g(\mathbf{x}^*) \geq \mathbf{0}$ .

Changing the sign of  $\mu^*$  as before, we obtain

1.  $\mu^* \leq \mathbf{0}$ .

2.  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} Dg(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\boldsymbol{\mu}^{*\top} g(\mathbf{x}^*) = 0$ .
4.  $h(\mathbf{x}^*) = \mathbf{0}$ .
5.  $g(\mathbf{x}^*) \geq \mathbf{0}$ .

For the problem

$$\begin{aligned} & \text{maximize} && f(\mathbf{x}) \\ & \text{subject to} && h(\mathbf{x}) = \mathbf{0} \\ & && g(\mathbf{x}) \geq \mathbf{0}, \end{aligned}$$

the KKT condition is exactly the same as in Theorem 21.1, except for the reversal of the inequality constraint.

**Example 21.3** In Figure 21.3, the two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are feasible points; that is,  $g(\mathbf{x}_1) \geq \mathbf{0}$  and  $g(\mathbf{x}_2) \geq \mathbf{0}$ , and they satisfy the KKT condition.

The point  $\mathbf{x}_1$  is a maximizer. The KKT condition for this point (with KKT multiplier  $\mu_1$ ) is

1.  $\mu_1 \geq 0$ .
2.  $\nabla f(\mathbf{x}_1) + \mu_1 \nabla g(\mathbf{x}_1) = \mathbf{0}$ .
3.  $\mu_1 g(\mathbf{x}_1) = 0$ .
4.  $g(\mathbf{x}_1) \geq \mathbf{0}$ .

The point  $\mathbf{x}_2$  is a minimizer of  $f$ . The KKT condition for this point (with KKT multiplier  $\mu_2$ ) is

1.  $\mu_2 \leq 0$ .
2.  $\nabla f(\mathbf{x}_2) + \mu_2 \nabla g(\mathbf{x}_2) = \mathbf{0}$ .
3.  $\mu_2 g(\mathbf{x}_2) = 0$ .
4.  $g(\mathbf{x}_2) \geq \mathbf{0}$ .

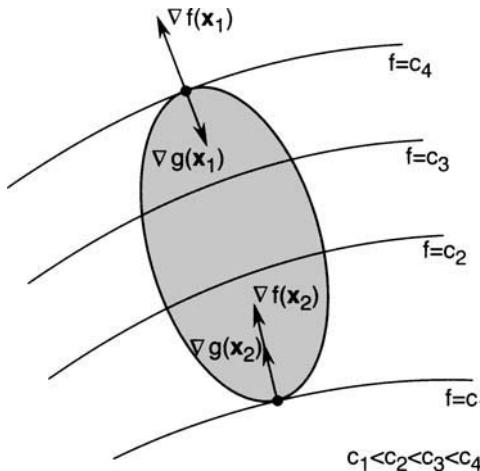
■

**Example 21.4** Consider the problem

$$\begin{aligned} & \text{minimize} && f(x_1, x_2) \\ & \text{subject to} && x_1, x_2 \geq 0, \end{aligned}$$

where

$$f(x_1, x_2) = x_1^2 + x_2^2 + x_1 x_2 - 3x_1.$$



**Figure 21.3** Points satisfying the KKT condition ( $x_1$  is a maximizer and  $x_2$  is a minimizer).

The KKT condition for this problem is

1.  $\mu = [\mu_1, \mu_2]^\top \leq \mathbf{0}$ .
2.  $Df(\mathbf{x}) + \mu^\top = \mathbf{0}^\top$ .
3.  $\mu^\top \mathbf{x} = 0$ .
4.  $\mathbf{x} \geq \mathbf{0}$ .

We have

$$Df(\mathbf{x}) = [2x_1 + x_2 - 3, x_1 + 2x_2].$$

This gives

$$\begin{aligned} 2x_1 + x_2 + \mu_1 &= 3, \\ x_1 + 2x_2 + \mu_2 &= 0, \\ \mu_1 x_1 + \mu_2 x_2 &= 0. \end{aligned}$$

We now have four variables, three equations, and the inequality constraints on each variable. To find a solution  $(\mathbf{x}^*, \mu^*)$ , we first try

$$\mu_1^* = 0, \quad x_2^* = 0,$$

which gives

$$x_1^* = \frac{3}{2}, \quad \mu_2^* = -\frac{3}{2}.$$

The above satisfies all the KKT and feasibility conditions. In a similar fashion, we can try

$$\mu_2^* = 0, \quad x_1^* = 0,$$

which gives

$$x_2^* = 0, \quad \mu_1^* = 3.$$

This point clearly violates the nonpositivity constraint on  $\mu_1^*$ .

The feasible point above satisfying the KKT condition is only a candidate for a minimizer. However, there is no guarantee that the point is indeed a minimizer, because the KKT condition is, in general, only necessary. A sufficient condition for a point to be a minimizer is given in the next section. ■

Example 21.4 is a special case of a more general problem of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

The KKT condition for this problem has the form

$$\begin{aligned} & \boldsymbol{\mu} \leq \mathbf{0}, \\ & \nabla f(\mathbf{x}) + \boldsymbol{\mu} = \mathbf{0}, \\ & \boldsymbol{\mu}^\top \mathbf{x} = 0, \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

From the above, we can eliminate  $\boldsymbol{\mu}$  to obtain

$$\begin{aligned} & \nabla f(\mathbf{x}) \geq \mathbf{0}, \\ & \mathbf{x}^\top \nabla f(\mathbf{x}) = 0, \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

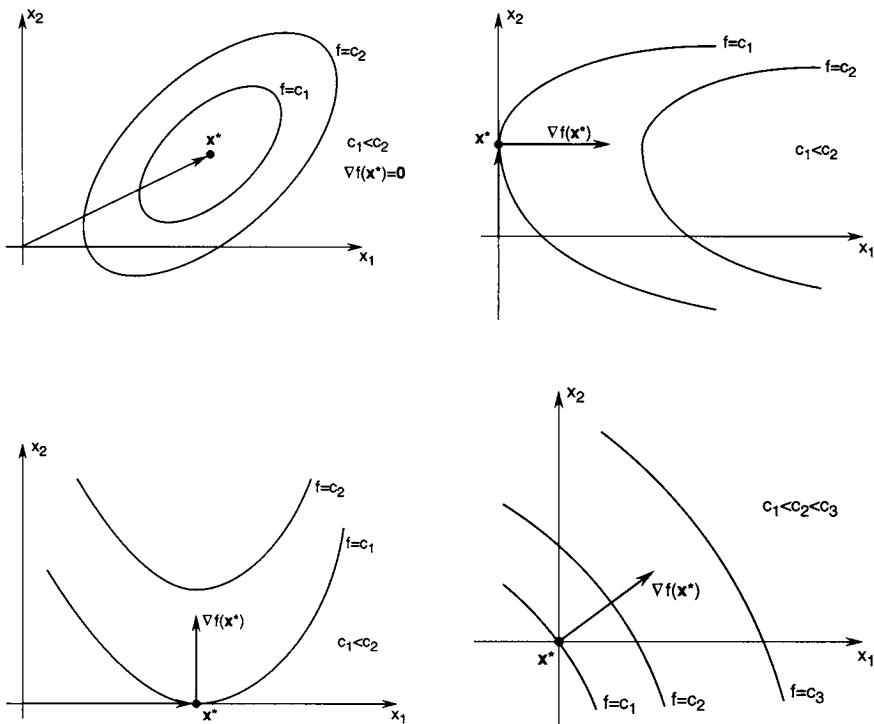
Some possible points in  $\mathbb{R}^2$  that satisfy these conditions are depicted in Figure 21.4.

For further results related to the KKT condition, we refer the reader to [90, Chapter 7].

## 21.2 Second-Order Conditions

As in the case of extremum problems with equality constraints, we can also give second-order necessary and sufficient conditions for extremum problems involving inequality constraints. For this, we need to define the following matrix:

$$\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{F}(\mathbf{x}) + [\boldsymbol{\lambda} \mathbf{H}(\mathbf{x})] + [\boldsymbol{\mu} \mathbf{G}(\mathbf{x})],$$



**Figure 21.4** Some possible points satisfying the KKT condition for problems with positive constraints. (Adapted from [13].)

where  $\mathbf{F}(\mathbf{x})$  is the Hessian matrix of  $f$  at  $\mathbf{x}$ , and the notation  $[\lambda \mathbf{H}(\mathbf{x})]$  represents

$$[\lambda \mathbf{H}(\mathbf{x})] = \lambda_1 \mathbf{H}_1(\mathbf{x}) + \cdots + \lambda_m \mathbf{H}_m(\mathbf{x}),$$

as before. Similarly, the notation  $[\mu \mathbf{G}(\mathbf{x})]$  represents

$$[\mu \mathbf{G}(\mathbf{x})] = \mu_1 \mathbf{G}_1(\mathbf{x}) + \cdots + \mu_p \mathbf{G}_p(\mathbf{x}),$$

where  $\mathbf{G}_k(\mathbf{x})$  is the Hessian of  $g_k$  at  $\mathbf{x}$ , given by

$$\mathbf{G}_k(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 g_k}{\partial x_1^2}(\mathbf{x}) & \cdots & \frac{\partial^2 g_k}{\partial x_n \partial x_1}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial^2 g_k}{\partial x_1 \partial x_n}(\mathbf{x}) & \cdots & \frac{\partial^2 g_k}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}.$$

In the following theorem, we use

$$T(\mathbf{x}^*) = \{\mathbf{y} \in \mathbb{R}^n : D\mathbf{h}(\mathbf{x}^*)\mathbf{y} = \mathbf{0}, Dg_j(\mathbf{x}^*)\mathbf{y} = 0, j \in J(\mathbf{x}^*)\},$$

that is, the tangent space to the surface defined by active constraints.

**Theorem 21.2 Second-Order Necessary Conditions.** *Let  $\mathbf{x}^*$  be a local minimizer of  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m \leq n$ ,  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ , and  $f, \mathbf{h}, \mathbf{g} \in C^2$ . Suppose that  $\mathbf{x}^*$  is regular. Then, there exist  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$  such that:*

1.  $\boldsymbol{\mu}^* \geq \mathbf{0}$ ,  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top$ ,  $\boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = 0$ .

2. For all  $\mathbf{y} \in T(\mathbf{x}^*)$  we have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \mathbf{y} \geq 0$ . □

*Proof.* Part 1 is simply a result of the KKT theorem. To prove part 2, we note that because the point  $\mathbf{x}^*$  is a local minimizer over  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ , it is also a local minimizer over  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, g_j(\mathbf{x}) = 0, j \in J(\mathbf{x}^*)\}$ ; that is, the point  $\mathbf{x}^*$  is a local minimizer with active constraints taken as equality constraints (see Exercise 21.16). Hence, the second-order necessary conditions for equality constraints (Theorem 20.4) are applicable here, which completes the proof. ■

We now state the second-order sufficient conditions for extremum problems involving inequality constraints. In the formulation of the result, we use the following set:

$$\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \{\mathbf{y} : D\mathbf{h}(\mathbf{x}^*) \mathbf{y} = \mathbf{0}, Dg_i(\mathbf{x}^*) \mathbf{y} = 0, i \in \tilde{J}(\mathbf{x}^*, \boldsymbol{\mu}^*)\},$$

where  $\tilde{J}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \{i : g_i(\mathbf{x}^*) = 0, \mu_i^* > 0\}$ . Note that  $\tilde{J}(\mathbf{x}^*, \boldsymbol{\mu}^*)$  is a subset of  $J(\mathbf{x}^*)$ :  $\tilde{J}(\mathbf{x}^*, \boldsymbol{\mu}^*) \subset J(\mathbf{x}^*)$ . This, in turn, implies that  $T(\mathbf{x}^*)$  is a subset of  $\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ :  $T(\mathbf{x}^*) \subset \tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ .

**Theorem 21.3 Second-Order Sufficient Conditions.** *Suppose that  $f, \mathbf{g}, \mathbf{h} \in C^2$  and there exist a feasible point  $\mathbf{x}^* \in \mathbb{R}^n$  and vectors  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$  such that:*

1.  $\boldsymbol{\mu}^* \geq \mathbf{0}$ ,  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top$ ,  $\boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = 0$ .

2. For all  $\mathbf{y} \in \tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ ,  $\mathbf{y} \neq \mathbf{0}$ , we have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \mathbf{y} > 0$ .

Then,  $\mathbf{x}^*$  is a strict local minimizer of  $f$  subject to  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ . □

*Proof.* For a proof of this theorem, we refer the reader to [88, p. 345]. ■

A result similar to Theorem 21.3 holds for a strict local maximizer, the only difference being that we need  $\boldsymbol{\mu}^* \leq \mathbf{0}$  and that  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  be negative definite on  $\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ .

**Example 21.5** Consider the following problem:

$$\begin{aligned} &\text{minimize} && x_1 x_2 \\ &\text{subject to} && x_1 + x_2 \geq 2 \\ & && x_2 \geq x_1. \end{aligned}$$

- a. Write down the KKT condition for this problem.
- b. Find all points (and KKT multipliers) satisfying the KKT condition. In each case, determine if the point is regular.
- c. Find all points in part b that also satisfy the SONC.
- d. Find all points in part c that also satisfy the SOSC.
- e. Find all points in part c that are local minimizers.

*Solution:*

- a. Write  $f(\mathbf{x}) = x_1 x_2$ ,  $g_1(\mathbf{x}) = 2 - x_1 - x_2$ , and  $g_2(\mathbf{x}) = x_1 - x_2$ . The KKT condition is

$$\begin{aligned}x_2 - \mu_1 + \mu_2 &= 0, \\x_1 - \mu_1 - \mu_2 &= 0, \\\mu_1(2 - x_1 - x_2) + \mu_2(x_1 - x_2) &= 0, \\\mu_1, \mu_2 &\geq 0, \\2 - x_1 - x_2 &\leq 0, \\x_1 - x_2 &\leq 0.\end{aligned}$$

- b. It is easy to check that  $\mu_1 \neq 0$  and  $\mu_2 \not> 0$ . This leaves us with only one solution to the KKT condition:  $x_1^* = x_2^* = 1$ ,  $\mu_1^* = 1$ ,  $\mu_2^* = 0$ . For this point, we have  $Dg_1(\mathbf{x}^*) = [-1, -1]$  and  $Dg_2(\mathbf{x}^*) = [1, -1]$ . Hence,  $\mathbf{x}^*$  is regular.
- c. Both constraints are active. Hence, because  $\mathbf{x}^*$  is regular,  $T(\mathbf{x}^*) = \{\mathbf{0}\}$ . This implies that the SONC is satisfied.

- d. Now,

$$\mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

Moreover,  $\tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*) = \{\mathbf{y} : [-1, -1]\mathbf{y} = 0\} = \{\mathbf{y} : y_1 = -y_2\}$ . Pick  $\mathbf{y} = [1, -1]^\top \in \tilde{T}(\mathbf{x}^*, \boldsymbol{\mu}^*)$ . We have  $\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) \mathbf{y} = -2 < 0$ , which means that the SOSC fails.

- e. In fact, the point  $\mathbf{x}^*$  is not a local minimizer. To see this, draw a picture of the constraint set and level sets of the objective function. Moving in the feasible direction  $[1, 1]^\top$ , the objective function increases; but moving in the feasible direction  $[-1, 1]^\top$ , the objective function decreases.

We now solve analytically the problem in Example 20.1 that we solved graphically earlier. ■

**Example 21.6** We wish to minimize  $f(\mathbf{x}) = (x_1 - 1)^2 + x_2 - 2$  subject to

$$\begin{aligned} h(\mathbf{x}) &= x_2 - x_1 - 1 = 0, \\ g(\mathbf{x}) &= x_1 + x_2 - 2 \leq 0. \end{aligned}$$

For all  $\mathbf{x} \in \mathbb{R}^2$ , we have

$$Dh(\mathbf{x}) = [-1, 1], \quad Dg(\mathbf{x}) = [1, 1].$$

Thus,  $\nabla h(\mathbf{x})$  and  $\nabla g(\mathbf{x})$  are linearly independent and hence all feasible points are regular. We first write the KKT condition. Because  $Df(\mathbf{x}) = [2x_1 - 2, 1]$ , we have

$$\begin{aligned} Df(\mathbf{x}) + \lambda Dh(\mathbf{x}) + \mu Dg(\mathbf{x}) &= [2x_1 - 2 - \lambda + \mu, 1 + \lambda + \mu] = \mathbf{0}^\top, \\ \mu(x_1 + x_2 - 2) &= 0, \\ \mu &\geq 0, \\ x_2 - x_1 - 1 &= 0, \\ x_1 + x_2 - 2 &\leq 0. \end{aligned}$$

To find points that satisfy the conditions above, we first try  $\mu > 0$ , which implies that  $x_1 + x_2 - 2 = 0$ . Thus, we are faced with a system of four linear equations

$$\begin{aligned} 2x_1 - 2 - \lambda + \mu &= 0, \\ 1 + \lambda + \mu &= 0, \\ x_2 - x_1 - 1 &= 0, \\ x_1 + x_2 - 2 &= 0. \end{aligned}$$

Solving the system of equations above, we obtain

$$x_1 = \frac{1}{2}, \quad x_2 = \frac{3}{2}, \quad \lambda = -1, \quad \mu = 0.$$

However, the above is not a legitimate solution to the KKT condition, because we obtained  $\mu = 0$ , which contradicts the assumption that  $\mu > 0$ .

In the second try, we assume that  $\mu = 0$ . Thus, we have to solve the system of equations

$$\begin{aligned} 2x_1 - 2 - \lambda &= 0, \\ 1 + \lambda &= 0, \\ x_2 - x_1 - 1 &= 0, \end{aligned}$$

and the solutions must satisfy

$$g(x_1, x_2) = x_1 + x_2 - 2 \leq 0.$$

Solving the equations above, we obtain

$$x_1 = \frac{1}{2}, \quad x_2 = \frac{3}{2}, \quad \lambda = -1.$$

Note that  $\mathbf{x}^* = [1/2, 3/2]^\top$  satisfies the constraint  $g(\mathbf{x}^*) \leq 0$ . The point  $\mathbf{x}^*$  satisfying the KKT necessary condition is therefore the candidate for being a minimizer.

We now verify if  $\mathbf{x}^* = [1/2, 3/2]^\top$ ,  $\lambda^* = -1$ ,  $\mu^* = 0$ , satisfy the second-order sufficient conditions. For this, we form the matrix

$$\begin{aligned} \mathbf{L}(\mathbf{x}^*, \lambda^*, \mu^*) &= \mathbf{F}(\mathbf{x}^*) + \lambda^* \mathbf{H}(\mathbf{x}^*) + \mu^* \mathbf{G}(\mathbf{x}^*) \\ &= \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} + (-1) \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + (0) \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}. \end{aligned}$$

We then find the subspace

$$\tilde{T}(\mathbf{x}^*, \mu^*) = \{\mathbf{y} : Dh(\mathbf{x}^*)\mathbf{y} = 0\}.$$

Note that because  $\mu^* = 0$ , the active constraint  $g(\mathbf{x}^*) = 0$  does not enter the computation of  $\tilde{T}(\mathbf{x}^*, \mu^*)$ . Note also that in this case,  $T(\mathbf{x}^*) = \{\mathbf{0}\}$ . We have

$$\tilde{T}(\mathbf{x}^*, \mu^*) = \{\mathbf{y} : [-1, 1]\mathbf{y} = 0\} = \{[a, a]^\top : a \in \mathbb{R}\}.$$

We then check for positive definiteness of  $\mathbf{L}(\mathbf{x}^*, \lambda^*, \mu^*)$  on  $\tilde{T}(\mathbf{x}^*, \mu^*)$ . We have

$$\mathbf{y}^\top \mathbf{L}(\mathbf{x}^*, \lambda^*, \mu^*) \mathbf{y} = [a, a] \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ a \end{bmatrix} = 2a^2.$$

Thus,  $\mathbf{L}(\mathbf{x}^*, \lambda^*, \mu^*)$  is positive definite on  $\tilde{T}(\mathbf{x}^*, \mu^*)$ . Observe that  $\mathbf{L}(\mathbf{x}^*, \lambda^*, \mu^*)$  is, in fact, only positive semidefinite on  $\mathbb{R}^2$ .

By the second-order sufficient conditions, we conclude that  $\mathbf{x}^* = [1/2, 3/2]^\top$  is a strict local minimizer. ■

## EXERCISES

**21.1** Consider the optimization problem

$$\begin{aligned} \text{minimize} \quad & x_1^2 + 4x_2^2 \\ \text{subject to} \quad & x_1^2 + 2x_2^2 \geq 4. \end{aligned}$$

- a. Find all the points that satisfy the KKT conditions.

- b. Apply the SOSC to determine the nature of the critical points from the previous part.

**21.2** Find local extremizers for:

- a.  $x_1^2 + x_2^2 - 2x_1 - 10x_2 + 26$  subject to  $\frac{1}{5}x_2 - x_1^2 \leq 0$ ,  $5x_1 + \frac{1}{2}x_2 \leq 5$ .
- b.  $x_1^2 + x_2^2$  subject to  $x_1 \geq 0$ ,  $x_2 \geq 0$ ,  $x_1 + x_2 \geq 5$ .
- c.  $x_1^2 + 6x_1x_2 - 4x_1 - 2x_2$  subject to  $x_1^2 + 2x_2 \leq 1$ ,  $2x_1 - 2x_2 \leq 1$ .

**21.3** Find local minimizers for  $x_1^2 + x_2^2$  subject to  $x_1^2 + 2x_1x_2 + x_2^2 = 1$ ,  $x_1^2 - x_2 \leq 0$ .

**21.4** Write down the Karush-Kuhn-Tucker condition for the optimization problem in Exercise 15.8.

**21.5** Consider the problem

$$\begin{aligned} & \text{minimize} && x_2 - (x_1 - 2)^3 + 3 \\ & \text{subject to} && x_2 \geq 1, \end{aligned}$$

where  $x_1$  and  $x_2$  are real variables. Answer each of the following questions, making sure that you give complete reasoning for your answers.

- a. Write down the KKT condition for the problem, and find all points that satisfy the condition. Check whether or not each point is regular.
- b. Determine whether or not the point(s) in part a satisfy the second-order necessary condition.
- c. Determine whether or not the point(s) in part b satisfy the second-order sufficient condition.

**21.6** Consider the problem

$$\begin{aligned} & \text{minimize} && x_2 \\ & \text{subject to} && x_2 \geq -(x_1 - 1)^2 + 3. \end{aligned}$$

- a. Find all points satisfying the KKT condition for the problem.
- b. For each point  $\mathbf{x}^*$  in part a, find  $T(\mathbf{x}^*)$ ,  $N(\mathbf{x}^*)$ , and  $\tilde{T}(\mathbf{x}^*)$ .
- c. Find the subset of points from part a that satisfy the second-order necessary condition.

**21.7** Consider the problem of optimizing (either minimizing or maximizing)  $(x_1 - 2)^2 + (x_2 - 1)^2$  subject to

$$\begin{aligned}x_2 - x_1^2 &\geq 0 \\2 - x_1 - x_2 &\geq 0 \\x_1 &\geq 0.\end{aligned}$$

The point  $\mathbf{x}^* = \mathbf{0}$  satisfies the KKT conditions.

- a. Does  $\mathbf{x}^*$  satisfy the FONC for minimization or maximization? What are the KKT multipliers?
- b. Does  $\mathbf{x}^*$  satisfy the SOSC? Carefully justify your answer.

**21.8** Consider the optimization problem

$$\begin{aligned}\text{minimize } & f(\mathbf{x}) \\ \text{subject to } & \mathbf{x} \in \Omega,\end{aligned}$$

where  $f(\mathbf{x}) = x_1 x_2^2$ , where  $\mathbf{x} = [x_1, x_2]^\top$ , and  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1 = x_2, x_1 \geq 0\}$ .

- a. Find all points satisfying the KKT condition.
- b. Do each of the points found in part a satisfy the second-order necessary condition?
- c. Do each of the points found in part a satisfy the second-order sufficient condition?

**21.9** Consider the problem

$$\begin{aligned}\text{minimize } & \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \\ \text{subject to } & x_1 + \cdots + x_n = 1 \\ & x_1, \dots, x_n \geq 0.\end{aligned}$$

- a. Write down the KKT condition for the problem.
- b. Define what it means for a feasible point  $\mathbf{x}^*$  to be *regular* in this particular problem. Are there any feasible points in this problem that are not regular? If yes, find them. If not, explain why not.

**21.10** Let  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $\mathbf{x}_0 \in \mathbb{R}^n$  be given, where  $g(\mathbf{x}_0) > 0$ . Consider the problem

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 \\ &\text{subject to} && g(\mathbf{x}) \leq 0. \end{aligned}$$

Suppose that  $\mathbf{x}^*$  is a solution to the problem and  $g \in \mathcal{C}^1$ . Use the KKT theorem to decide which of the following equations/inequalities hold:

- i.  $g(\mathbf{x}^*) < 0$ .
- ii.  $g(\mathbf{x}^*) = 0$ .
- iii.  $(\mathbf{x}^* - \mathbf{x}_0)^\top \nabla g(\mathbf{x}^*) < 0$ .
- iv.  $(\mathbf{x}^* - \mathbf{x}_0)^\top \nabla g(\mathbf{x}^*) = 0$ .
- v.  $(\mathbf{x}^* - \mathbf{x}_0)^\top \nabla g(\mathbf{x}^*) > 0$ .

**21.11** Consider a square room with corners located at  $[0, 0]^\top$ ,  $[0, 2]^\top$ ,  $[2, 0]^\top$ , and  $[2, 2]^\top$  (in  $\mathbb{R}^2$ ). We wish to find the point in the room that is closest to the point  $[3, 4]^\top$ .

- a. Guess which point in the room is the closest point in the room to the point  $[3, 4]^\top$ .
- b. Use the second-order sufficient conditions to prove that the point you have guessed is a strict local minimizer.

*Hint:* Minimizing the distance is the same as minimizing the square distance.

**21.12** Consider the *quadratic programming problem*

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \\ &\text{subject to} && \mathbf{A} \mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and  $\mathbf{b} \geq \mathbf{0}$ . Find all points satisfying the KKT condition.

**21.13** Consider the linear programming problem

$$\begin{aligned} &\text{minimize} && ax_1 + bx_2 \\ &\text{subject to} && cx_1 + dx_2 = e \\ &&& x_1, x_2 \geq 0, \end{aligned}$$

where  $a, b, c, d, e \in \mathbb{R}$  are all nonzero constants. Suppose that  $\mathbf{x}^*$  is an optimal basic feasible solution to the problem.

- a. Write down the Karush-Kuhn-Tucker condition involving  $\mathbf{x}^*$  (specifying clearly the number of Lagrange and KKT multipliers).
- b. Is  $\mathbf{x}^*$  regular? Explain.
- c. Find the tangent space  $T(\mathbf{x}^*)$  (defined by the active constraints) for this problem.
- d. Assume that the relative cost coefficients of all nonbasic variables are strictly positive. Does  $\mathbf{x}^*$  satisfy the second-order sufficient condition? Explain.

**21.14** Consider the problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{0}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ , is of full rank. Use the KKT theorem to show that if there exists a solution, then the optimal objective function value is 0.

**21.15** Consider a linear programming problem in standard form (see Chapter 15).

- a. Write down the Karush-Kuhn-Tucker condition for the problem.
- b. Use part a to show that if there exists an optimal feasible solution to the linear program, then there exists a feasible solution to the corresponding dual problem that achieves an objective function value that is the same as the optimal value of the primal (compare this with Theorem 17.1).
- c. Use parts a and b to prove that if  $\mathbf{x}^*$  is an optimal feasible solutions of the primal, then there exists a feasible solution  $\boldsymbol{\lambda}^*$  to the dual such that  $(\mathbf{c}^\top - \boldsymbol{\lambda}^{*\top} \mathbf{A})\mathbf{x}^* = 0$  (compare this with Theorem 17.3).

**21.16** Consider the constraint set  $S = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$ . Let  $\mathbf{x}^* \in S$  be a regular local minimizer of  $f$  over  $S$  and  $J(\mathbf{x}^*)$  the index set of active inequality constraints. Show that  $\mathbf{x}^*$  is also a regular local minimizer of  $f$  over the set  $S' = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, g_j(\mathbf{x}) = 0, j \in J(\mathbf{x}^*)\}$ .

**21.17** Solve the following optimization problem using the second-order sufficient conditions:

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && x_1^2 - x_2 - 4 \leq 0 \\ & && x_2 - x_1 - 2 \leq 0. \end{aligned}$$

See Figure 22.1 for a graphical illustration of the problem.

**21.18** Solve the following optimization problem using the second-order sufficient conditions:

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && x_1 - x_2^2 - 4 \geq 0 \\ & && x_1 - 10 \leq 0. \end{aligned}$$

See Figure 22.2 for a graphical illustration of the problem.

**21.19** Consider the problem

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && 4 - x_1 - x_2^2 \leq 0 \\ & && 3x_2 - x_1 \leq 0 \\ & && -3x_2 - x_1 \leq 0. \end{aligned}$$

Figure 22.3 gives a graphical illustration of the problem. Deduce from the figure that the problem has two strict local minimizers, and use the second-order sufficient conditions to verify the graphical solutions.

**21.20** Consider the following optimization problem with an inequality constraint:

$$\begin{aligned} & \text{minimize} && 3x_1 \\ & \text{subject to} && x_1 + x_2^2 \geq 2. \end{aligned}$$

- a. Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the KKT (first-order necessary) condition?
- b. Does the point  $\mathbf{x}^* = [2, 0]^\top$  satisfy the second-order necessary condition (for problems with inequality constraints)?
- c. Is the point  $\mathbf{x}^* = [2, 0]^\top$  a local minimizer?

(See Exercise 6.15 for a similar problem treated using set-constrained methods.)

**21.21** Consider the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{x}\|^2 \\ & \text{subject to} && \mathbf{a}^\top \mathbf{x} = b \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{a} \in \mathbb{R}^n$ ,  $\mathbf{a} \geq \mathbf{0}$ , and  $b \in \mathbb{R}$ ,  $b > 0$ . Show that if a solution to the problem exists, then it is unique, and find an expression for it in terms of  $\mathbf{a}$  and  $b$ .

**21.22** Consider the problem

$$\begin{aligned} &\text{minimize} && (x_1 - a)^2 + (x_2 - b)^2, \quad x_1, x_2 \in \mathbb{R} \\ &\text{subject to} && x_1^2 + x_2^2 \leq 1, \end{aligned}$$

where  $a, b \in \mathbb{R}$  are given constants satisfying  $a^2 + b^2 \geq 1$ .

- a. Let  $\mathbf{x}^* = [x_1^*, x_2^*]^\top$  be a solution to the problem. Use the first-order necessary conditions for unconstrained optimization to show that  $(x_1^*)^2 + (x_2^*)^2 = 1$ .
- b. Use the KKT theorem to show that the solution  $\mathbf{x}^* = [x_1^*, x_2^*]^\top$  is unique and has the form  $x_1^* = \alpha a$ ,  $x_2^* = \alpha b$ , where  $\alpha \in \mathbb{R}$  is a positive constant.
- c. Find an expression for  $\alpha$  (from part b) in terms of  $a$  and  $b$ .

**21.23** Consider the problem

$$\begin{aligned} &\text{minimize} && x_1^2 + (x_2 + 1)^2, \quad x_1, x_2 \in \mathbb{R} \\ &\text{subject to} && x_2 \geq \exp(x_1) \end{aligned}$$

[ $\exp(x) = e^x$  is the exponential of  $x$ ]. Let  $\mathbf{x}^* = [x_1^*, x_2^*]^\top$  be the solution to the problem.

- a. Write down the KKT condition that must be satisfied by  $\mathbf{x}^*$ .
- b. Prove that  $x_2^* = \exp(x_1^*)$ .
- c. Prove that  $-2 < x_1^* < 0$ .

**21.24** Consider the problem

$$\begin{aligned} &\text{minimize} && \mathbf{c}^\top \mathbf{x} + 8 \\ &\text{subject to} && \frac{1}{2} \|\mathbf{x}\|^2 \leq 1, \end{aligned}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{c} \neq \mathbf{0}$ . Suppose that  $\mathbf{x}^* = \alpha \mathbf{e}$  is a solution to the problem, where  $\alpha \in \mathbb{R}$  and  $\mathbf{e} = [1, \dots, 1]^\top$ , and the corresponding objective value is 4.

- a. Show that  $\|\mathbf{x}^*\|^2 = 2$ .
- b. Find  $\alpha$  and  $\mathbf{c}$  (they may depend on  $n$ ).

**21.25** Consider the problem with equality constraint

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}. \end{aligned}$$

We can convert the above into the equivalent optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \frac{1}{2} \|\mathbf{h}(\mathbf{x})\|^2 \leq 0. \end{aligned}$$

Write down the KKT condition for the equivalent problem (with inequality constraint) and explain why the KKT theorem cannot be applied in this case.

## CHAPTER 22

---

# CONVEX OPTIMIZATION PROBLEMS

---

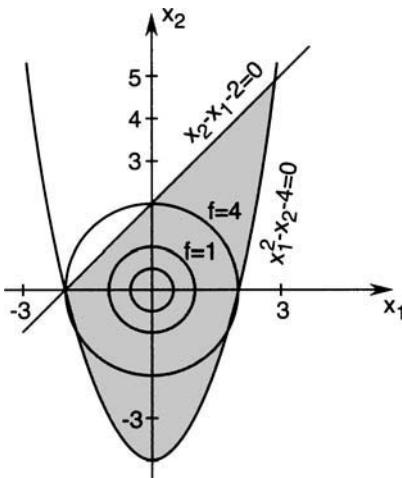
### 22.1 Introduction

The optimization problems posed at the beginning of this part are, in general, very difficult to solve. The source of these difficulties may be in the objective function or the constraints. Even if the objective function is simple and “well-behaved,” the nature of the constraints may make the problem difficult to solve. We illustrate some of these difficulties in the following examples.

**Example 22.1** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && x_2 - x_1 - 2 \leq 0 \\ & && x_1^2 - x_2 - 4 \leq 0. \end{aligned}$$

The problem is depicted in Figure 22.1, where, as we can see, the constrained minimizer is the same as the unconstrained minimizer. At the minimizer, all the constraints are inactive. If we had only known this fact, we could have approached this problem as an unconstrained optimization problem using techniques from Part II. ■



**Figure 22.1** Situation where the constrained and the unconstrained minimizers are the same.

**Example 22.2** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && x_1 - 10 \leq 0 \\ & && x_1 - x_2 - 4 \geq 0. \end{aligned}$$

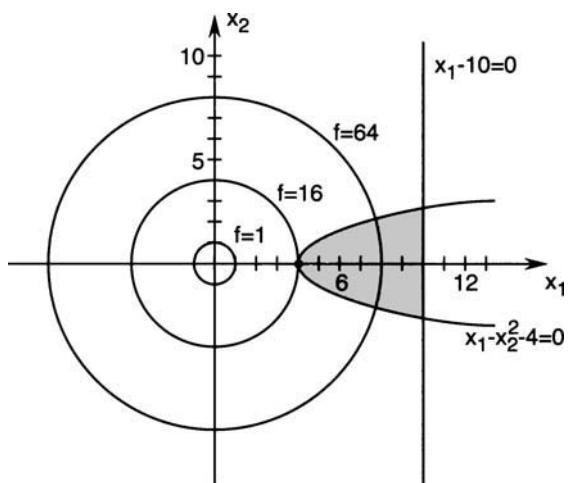
The problem is depicted in Figure 22.2. At the solution, only one constraint is active. If we had only known about this we could have handled this problem as a constrained optimization problem using the Lagrange multiplier method. ■

**Example 22.3** Consider the optimization problem

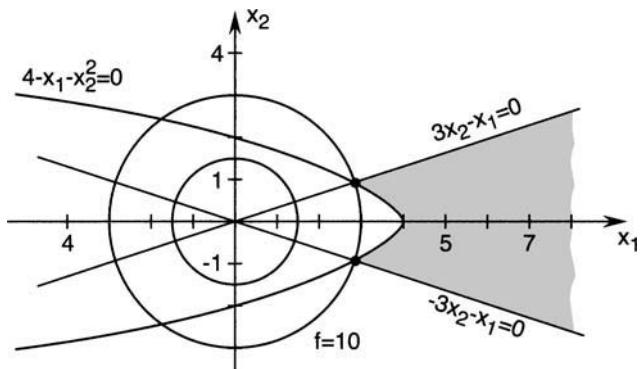
$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 \\ & \text{subject to} && 4 - x_1 - x_2^2 \leq 0 \\ & && 3x_2 - x_1 \leq 0 \\ & && -3x_2 - x_1 \leq 0. \end{aligned}$$

The problem is depicted in Figure 22.3. This example illustrates the situation where the constraints introduce local minimizers, even though the objective function itself has only one unconstrained global minimizer. ■

Some of the difficulties illustrated in the examples above can be eliminated if we restrict our problems to convex feasible regions. Admittedly, some important real-life problems do not fit into this framework. On the other hand,



**Figure 22.2** Situation where only one constraint is active.



**Figure 22.3** Situation where the constraints introduce local minimizers.

it is possible to give results of a *global* nature for this class of optimization problems. In the next section, we introduce the notion of a *convex function*, which plays an important role in our subsequent treatment of such problems.

## 22.2 Convex Functions

We begin with a definition of the graph of a real-valued function.

**Definition 22.1** The *graph* of  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^n$ , is the set of points in  $\Omega \times \mathbb{R} \subset \mathbb{R}^{n+1}$  given by

$$\left\{ \begin{bmatrix} \mathbf{x} \\ f(\mathbf{x}) \end{bmatrix} : \mathbf{x} \in \Omega \right\}.$$

■

We can visualize the graph of  $f$  as simply the set of points on a “plot” of  $f(\mathbf{x})$  versus  $\mathbf{x}$  (see Figure 22.4). We next define the epigraph of a real-valued function.

**Definition 22.2** The *epigraph* of a function  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^n$ , denoted  $\text{epi}(f)$ , is the set of points in  $\Omega \times \mathbb{R}$  given by

$$\text{epi}(f) = \left\{ \begin{bmatrix} \mathbf{x} \\ \beta \end{bmatrix} : \mathbf{x} \in \Omega, \beta \in \mathbb{R}, \beta \geq f(\mathbf{x}) \right\}.$$

■

The epigraph  $\text{epi}(f)$  of a function  $f$  is simply the set of points in  $\Omega \times \mathbb{R}$  on or above the graph of  $f$  (see Figure 22.4). We can also think of  $\text{epi}(f)$  as a subset of  $\mathbb{R}^{n+1}$ .

Recall that a set  $\Omega \subset \mathbb{R}^n$  is convex if for every  $\mathbf{x}_1, \mathbf{x}_2 \in \Omega$  and  $\alpha \in (0, 1)$ ,  $\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2 \in \Omega$  (see Section 4.3). We now introduce the notion of a convex function.

**Definition 22.3** A function  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^n$ , is *convex* on  $\Omega$  if its epigraph is a convex set. ■

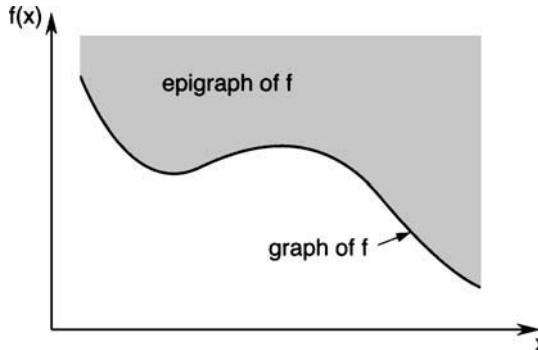
**Theorem 22.1** If a function  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^n$ , is convex on  $\Omega$ , then  $\Omega$  is a convex set. □

*Proof.* We prove this theorem by contraposition. Suppose that  $\Omega$  is not a convex set. Then, there exist two points  $\mathbf{y}_1$  and  $\mathbf{y}_2$  such that for some  $\alpha \in (0, 1)$ ,

$$\mathbf{z} = \alpha\mathbf{y}_1 + (1 - \alpha)\mathbf{y}_2 \notin \Omega.$$

Let

$$\beta_1 = f(\mathbf{y}_1), \quad \beta_2 = f(\mathbf{y}_2).$$



**Figure 22.4** Graph and epigraph of a function  $f : \mathbb{R} \rightarrow \mathbb{R}$ .

Then, the pairs

$$\begin{bmatrix} \mathbf{y}_1 \\ \beta_1 \end{bmatrix}, \begin{bmatrix} \mathbf{y}_2 \\ \beta_2 \end{bmatrix}$$

belong to the graph of  $f$ , and hence also the epigraph of  $f$ . Let

$$\mathbf{w} = \alpha \begin{bmatrix} \mathbf{y}_1 \\ \beta_1 \end{bmatrix} + (1 - \alpha) \begin{bmatrix} \mathbf{y}_2 \\ \beta_2 \end{bmatrix}.$$

We have

$$\mathbf{w} = \begin{bmatrix} \mathbf{z} \\ \alpha\beta_1 + (1 - \alpha)\beta_2 \end{bmatrix}.$$

But note that  $\mathbf{w} \notin \text{epi}(f)$ , because  $\mathbf{z} \notin \Omega$ . Therefore,  $\text{epi}(f)$  is not convex, and hence  $f$  is not a convex function. ■

The next theorem gives a very useful characterization of convex functions. This characterization is often used as a definition for a convex function.

**Theorem 22.2** *A function  $f : \Omega \rightarrow \mathbb{R}$  defined on a convex set  $\Omega \subset \mathbb{R}^n$  is convex if and only if for all  $\mathbf{x}, \mathbf{y} \in \Omega$  and all  $\alpha \in (0, 1)$ , we have*

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}).$$

□

*Proof.*  $\Leftarrow$ : Assume that for all  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ ,

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}).$$

Let  $[\mathbf{x}^\top, a]^\top$  and  $[\mathbf{y}^\top, b]^\top$  be two points in  $\text{epi}(f)$ , where  $a, b \in \mathbb{R}$ . From the definition of  $\text{epi}(f)$  it follows that

$$f(\mathbf{x}) \leq a, \quad f(\mathbf{y}) \leq b.$$

Therefore, using the first inequality above, we have

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha a + (1 - \alpha)b.$$

Because  $\Omega$  is convex,  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \Omega$ . Hence,

$$\begin{bmatrix} \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \\ \alpha a + (1 - \alpha)b \end{bmatrix} \in \text{epi}(f),$$

which implies that  $\text{epi}(f)$  is a convex set, and hence  $f$  is a convex function.

$\Rightarrow$ : Assume that  $f : \Omega \rightarrow \mathbb{R}$  is a convex function. Let  $\mathbf{x}, \mathbf{y} \in \Omega$  and

$$f(\mathbf{x}) = a, \quad f(\mathbf{y}) = b.$$

Thus,

$$\begin{bmatrix} \mathbf{x} \\ a \end{bmatrix}, \begin{bmatrix} \mathbf{y} \\ b \end{bmatrix} \in \text{epi}(f).$$

Because  $f$  is a convex function, its epigraph is a convex subset of  $\mathbb{R}^{n+1}$ . Therefore, for all  $\alpha \in (0, 1)$ , we have

$$\alpha \begin{bmatrix} \mathbf{x} \\ a \end{bmatrix} + (1 - \alpha) \begin{bmatrix} \mathbf{y} \\ b \end{bmatrix} = \begin{bmatrix} \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \\ \alpha a + (1 - \alpha)b \end{bmatrix} \in \text{epi}(f).$$

The above implies that for all  $\alpha \in (0, 1)$ ,

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha a + (1 - \alpha)b = \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}).$$

This completes the proof. ■

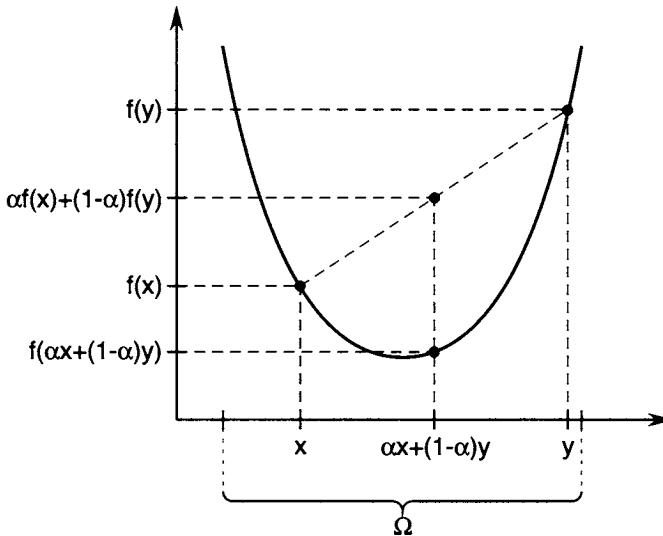
A geometric interpretation of Theorem 22.2 is given in Figure 22.5. The theorem states that if  $f : \Omega \rightarrow \mathbb{R}$  is a convex function over a convex set  $\Omega$ , then for all  $\mathbf{x}, \mathbf{y} \in \Omega$ , the points on the line segment in  $\mathbb{R}^{n+1}$  connecting  $[\mathbf{x}^\top, f(\mathbf{x})]^\top$  and  $[\mathbf{y}^\top, f(\mathbf{y})]^\top$  must lie on or above the graph of  $f$ .

Using Theorem 22.2, it is straightforward to show that any nonnegative scaling of a convex function is convex, and that the sum of convex functions is convex.

**Theorem 22.3** Suppose that  $f$ ,  $f_1$ , and  $f_2$  are convex functions. Then, for any  $a \geq 0$ , the function  $af$  is convex. Moreover,  $f_1 + f_2$  is convex. □

*Proof.* Let  $\mathbf{x}, \mathbf{y} \in \Omega$  and  $\alpha \in (0, 1)$ . Fix  $a \geq 0$ . For convenience, write  $\bar{f} = af$ . We have

$$\begin{aligned} \bar{f}(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) &= af(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \\ &\leq a(\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})) \text{ because } f \text{ is convex and } a \geq 0 \\ &= \alpha(af(\mathbf{x})) + (1 - \alpha)(af(\mathbf{y})) \\ &= \alpha\bar{f}(\mathbf{x}) + (1 - \alpha)\bar{f}(\mathbf{y}), \end{aligned}$$



**Figure 22.5** Geometric interpretation of Theorem 22.2.

which implies that  $\bar{f}$  is convex.

Next, write  $f_3 = f_1 + f_2$ . We have

$$\begin{aligned}
 f_3(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) &= f_1(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) + f_2(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \\
 &\leq (\alpha f_1(\mathbf{x}) + (1 - \alpha)f_1(\mathbf{y})) + (\alpha f_2(\mathbf{x}) + (1 - \alpha)f_2(\mathbf{y})) \\
 &\quad \text{by convexity of } f_1 \text{ and } f_2 \\
 &= \alpha(f_1(\mathbf{x}) + f_2(\mathbf{x})) + (1 - \alpha)(f_1(\mathbf{y}) + f_2(\mathbf{y})) \\
 &= \alpha f_3(\mathbf{x}) + (1 - \alpha)f_3(\mathbf{y}),
 \end{aligned}$$

which implies that  $f_3$  is convex. ■

Theorem 22.3 implies that for any given collection of convex functions  $f_1, \dots, f_\ell$  and nonnegative numbers  $c_1, \dots, c_\ell$ , the function  $c_1 f_1 + \dots + c_\ell f_\ell$  is convex. Using a method of proof similar to that used in Theorem 22.3, it is similarly straightforward to show that the function  $\max\{f_1, \dots, f_\ell\}$  is convex (see Exercise 22.6).

We now define the notion of strict convexity.

**Definition 22.4** A function  $f : \Omega \rightarrow \mathbb{R}$  on a convex set  $\Omega \subset \mathbb{R}^n$  is *strictly convex* if for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,  $\mathbf{x} \neq \mathbf{y}$ , and  $\alpha \in (0, 1)$ , we have

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) < \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}).$$
■

From this definition, we see that for a strictly convex function, all points on the open line segment connecting the points  $[\mathbf{x}^\top, f(\mathbf{x})]^\top$  and  $[\mathbf{y}^\top, f(\mathbf{y})]^\top$  lie (strictly) above the graph of  $f$ .

**Definition 22.5** A function  $f : \Omega \rightarrow \mathbb{R}$  on a convex set  $\Omega \subset \mathbb{R}^n$  is (strictly) *concave* if  $-f$  is (strictly) convex. ■

Note that the graph of a strictly concave function always lies above the line segment connecting any two points on its graph.

To show that a function is not convex, we need only produce a pair of points  $\mathbf{x}, \mathbf{y} \in \Omega$  and an  $\alpha \in (0, 1)$  such that the inequality in Theorem 22.2 is violated.

**Example 22.4** Let  $f(\mathbf{x}) = x_1 x_2$ . Is  $f$  convex over  $\Omega = \{\mathbf{x} : x_1 \geq 0, x_2 \geq 0\}$ ?

The answer is no. Take, for example,  $\mathbf{x} = [1, 2]^\top \in \Omega$  and  $\mathbf{y} = [2, 1]^\top \in \Omega$ . Then,

$$\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} = \begin{bmatrix} 2 - \alpha \\ 1 + \alpha \end{bmatrix}.$$

Hence,

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) = (2 - \alpha)(1 + \alpha) = 2 + \alpha - \alpha^2$$

and

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) = 2.$$

If, for example,  $\alpha = 1/2 \in (0, 1)$ , then

$$f\left(\frac{1}{2}\mathbf{x} + \frac{1}{2}\mathbf{y}\right) = \frac{9}{4} > \frac{1}{2}f(\mathbf{x}) + \frac{1}{2}f(\mathbf{y}),$$

which shows that  $f$  is not convex over  $\Omega$ . ■

Example 22.4 is an illustration of the following general result.

**Proposition 22.1** A quadratic form  $f : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^n$ , given by  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ ,  $\mathbf{Q} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{Q} = \mathbf{Q}^\top$ , is convex on  $\Omega$  if and only if for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,  $(\mathbf{x} - \mathbf{y})^\top \mathbf{Q} (\mathbf{x} - \mathbf{y}) \geq 0$ . □

*Proof.* The result follows from Theorem 22.2. Indeed, the function  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q} \mathbf{x}$  is convex if and only if for every  $\alpha \in (0, 1)$ , and every  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , we have

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}),$$

or, equivalently,

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) - f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \geq 0.$$

Substituting for  $f$  into the left-hand side of this equation yields

$$\begin{aligned} \alpha\mathbf{x}^\top Q\mathbf{x} + (1-\alpha)\mathbf{y}^\top Q\mathbf{y} - (\alpha\mathbf{x} + (1-\alpha)\mathbf{y})^\top Q(\alpha\mathbf{x} + (1-\alpha)\mathbf{y}) \\ = \alpha\mathbf{x}^\top Q\mathbf{x} + \mathbf{y}^\top Q\mathbf{y} - \alpha\mathbf{y}^\top Q\mathbf{y} - \alpha^2\mathbf{x}^\top Q\mathbf{x} \\ - (2\alpha - 2\alpha^2)\mathbf{x}^\top Q\mathbf{y} - (1 - 2\alpha + \alpha^2)\mathbf{y}^\top Q\mathbf{y} \\ = \alpha(1-\alpha)\mathbf{x}^\top Q\mathbf{x} - 2\alpha(1-\alpha)\mathbf{x}^\top Q\mathbf{y} + \alpha(1-\alpha)\mathbf{y}^\top Q\mathbf{y} \\ = \alpha(1-\alpha)(\mathbf{x} - \mathbf{y})^\top Q(\mathbf{x} - \mathbf{y}). \end{aligned}$$

Therefore,  $f$  is convex if and only if

$$\alpha(1-\alpha)(\mathbf{x} - \mathbf{y})^\top Q(\mathbf{x} - \mathbf{y}) \geq 0,$$

which proves the result. ■

**Example 22.5** In Example 22.4,  $f(\mathbf{x}) = x_1x_2$ , which can be written as  $f(\mathbf{x}) = \mathbf{x}^\top Q\mathbf{x}$ , where

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

Let  $\Omega = \{\mathbf{x} : \mathbf{x} \geq \mathbf{0}\}$ , and  $\mathbf{x} = [2, 2]^\top \in \Omega$ ,  $\mathbf{y} = [1, 3]^\top \in \Omega$ . We have

$$\mathbf{y} - \mathbf{x} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

and

$$(\mathbf{y} - \mathbf{x})^\top Q(\mathbf{y} - \mathbf{x}) = \frac{1}{2}[-1, 1] \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = -1 < 0.$$

Hence, by Proposition 22.1,  $f$  is not convex on  $\Omega$ . ■

Differentiable convex functions can be characterized using the following theorem.

**Theorem 22.4** Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , be defined on an open convex set  $\Omega \subset \mathbb{R}^n$ . Then,  $f$  is convex on  $\Omega$  if and only if for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

□

*Proof.*  $\Rightarrow$ : Suppose that  $f : \Omega \rightarrow \mathbb{R}$  is differentiable and convex. Then, by Theorem 22.2, for any  $\mathbf{y}, \mathbf{x} \in \Omega$  and  $\alpha \in (0, 1)$  we have

$$f(\alpha\mathbf{y} + (1-\alpha)\mathbf{x}) \leq \alpha f(\mathbf{y}) + (1-\alpha)f(\mathbf{x}).$$

Rearranging terms yields

$$f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x}) \leq \alpha(f(\mathbf{y}) - f(\mathbf{x})).$$

Upon dividing both sides of this inequality by  $\alpha$ , we get

$$\frac{f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\alpha} \leq f(\mathbf{y}) - f(\mathbf{x}).$$

If we now take the limit as  $\alpha \rightarrow 0$  and apply the definition of the directional derivative of  $f$  at  $\mathbf{x}$  in the direction  $\mathbf{y} - \mathbf{x}$  (see Section 6.2), we get

$$Df(\mathbf{x})(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$$

or

$$f(\mathbf{y}) \geq f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

$\Leftarrow$ : Assume that  $\Omega$  is convex,  $f : \Omega \rightarrow \mathbb{R}$  is differentiable, and for all  $\mathbf{x}, \mathbf{y} \in \Omega$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

Let  $\mathbf{u}, \mathbf{v} \in \Omega$  and  $\alpha \in (0, 1)$ . Because  $\Omega$  is convex,

$$\mathbf{w} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v} \in \Omega.$$

We also have

$$f(\mathbf{u}) \geq f(\mathbf{w}) + Df(\mathbf{w})(\mathbf{u} - \mathbf{w})$$

and

$$f(\mathbf{v}) \geq f(\mathbf{w}) + Df(\mathbf{w})(\mathbf{v} - \mathbf{w}).$$

Multiplying the first of these inequalities by  $\alpha$  and the second by  $(1 - \alpha)$  and then adding them together yields

$$\alpha f(\mathbf{u}) + (1 - \alpha)f(\mathbf{v}) \geq f(\mathbf{w}) + Df(\mathbf{w})(\alpha\mathbf{u} + (1 - \alpha)\mathbf{v} - \mathbf{w}).$$

But

$$\mathbf{w} = \alpha\mathbf{u} + (1 - \alpha)\mathbf{v}.$$

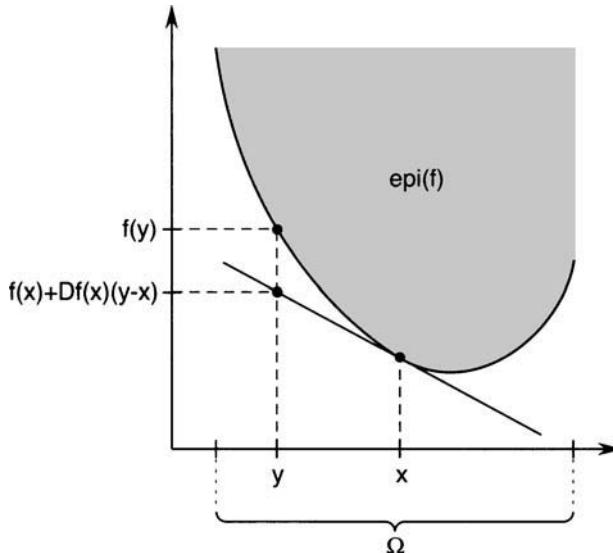
Hence,

$$\alpha f(\mathbf{u}) + (1 - \alpha)f(\mathbf{v}) \geq f(\alpha\mathbf{u} + (1 - \alpha)\mathbf{v}).$$

Hence, by Theorem 22.2,  $f$  is a convex function. ■

In Theorem 22.4, the assumption that  $\Omega$  be open is not necessary, as long as  $f \in \mathcal{C}^1$  on some open set that contains  $\Omega$  (e.g.,  $f \in \mathcal{C}^1$  on  $\mathbb{R}^n$ ).

A geometric interpretation of Theorem 22.4 is given in Figure 22.6. To explain the interpretation, let  $\mathbf{x}_0 \in \Omega$ . The function  $\ell(\mathbf{x}) = f(\mathbf{x}_0) + Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)$  is the linear approximation to  $f$  at  $\mathbf{x}_0$ . The theorem says that the graph of  $f$  always lies above its linear approximation at any point. In other words,



**Figure 22.6** Geometric interpretation of Theorem 22.4.

the linear approximation to a convex function  $f$  at any point of its domain lies below  $\text{epi}(f)$ .

This geometric idea leads to a generalization of the gradient to the case where  $f$  is not differentiable. Let  $f : \Omega \rightarrow \mathbb{R}$  be defined on an open convex set  $\Omega \subset \mathbb{R}^n$ . A vector  $\mathbf{g} \in \mathbb{R}^n$  is said to be a *subgradient* of  $f$  at a point  $\mathbf{x} \in \Omega$  if for all  $\mathbf{y} \in \Omega$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}).$$

As in the case of the standard gradient, if  $\mathbf{g}$  is a subgradient, then for a given  $\mathbf{x}_0 \in \Omega$ , the function  $\ell(\mathbf{x}) = f(\mathbf{x}_0) + \mathbf{g}^\top (\mathbf{x} - \mathbf{x}_0)$  lies below  $\text{epi}(f)$ .

For functions that are twice continuously differentiable, the following theorem gives another possible characterization of convexity.

**Theorem 22.5** *Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $f \in C^2$ , be defined on an open convex set  $\Omega \subset \mathbb{R}^n$ . Then,  $f$  is convex on  $\Omega$  if and only if for each  $\mathbf{x} \in \Omega$ , the Hessian  $\mathbf{F}(\mathbf{x})$  of  $f$  at  $\mathbf{x}$  is a positive semidefinite matrix.*  $\square$

*Proof.*  $\Leftarrow$ : Let  $\mathbf{x}, \mathbf{y} \in \Omega$ . Because  $f \in C^2$ , by Taylor's theorem there exists  $\alpha \in (0, 1)$  such that

$$f(\mathbf{y}) = f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top \mathbf{F}(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}).$$

Because  $\mathbf{F}(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))$  is positive semidefinite,

$$(\mathbf{y} - \mathbf{x})^\top \mathbf{F}(\alpha\mathbf{y} + (1 - \alpha)\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0.$$

Therefore, we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}),$$

which implies that  $f$  is convex, by Theorem 22.4.

$\Rightarrow$ : We use contraposition. Assume that there exists  $\mathbf{x} \in \Omega$  such that  $\mathbf{F}(\mathbf{x})$  is not positive semidefinite. Therefore, there exists  $\mathbf{d} \in \mathbb{R}^n$  such that  $\mathbf{d}^\top \mathbf{F}(\mathbf{x}) \mathbf{d} < 0$ . By assumption,  $\Omega$  is open; thus, the point  $\mathbf{x}$  is an interior point. By the continuity of the Hessian matrix, there exists a nonzero  $s \in \mathbb{R}$  such that  $\mathbf{x} + s\mathbf{d} \in \Omega$ , and if we write  $\mathbf{y} = \mathbf{x} + s\mathbf{d}$ , then for all points  $\mathbf{z}$  on the line segment joining  $\mathbf{x}$  and  $\mathbf{y}$ , we have  $\mathbf{d}^\top \mathbf{F}(\mathbf{z}) \mathbf{d} < 0$ . By Taylor's theorem there exists  $\alpha \in (0, 1)$  such that

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top \mathbf{F}(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}) \\ &= f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}s^2 \mathbf{d}^\top \mathbf{F}(\mathbf{x} + \alpha s\mathbf{d}) \mathbf{d}. \end{aligned}$$

Because  $\alpha \in (0, 1)$ , the point  $\mathbf{x} + \alpha s\mathbf{d}$  is on the line segment joining  $\mathbf{x}$  and  $\mathbf{y}$ , and therefore

$$\mathbf{d}^\top \mathbf{F}(\mathbf{x} + \alpha s\mathbf{d}) \mathbf{d} < 0.$$

Because  $s \neq 0$ , we have  $s^2 > 0$ , and hence

$$f(\mathbf{y}) < f(\mathbf{x}) + Df(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

Therefore, by Theorem 22.4,  $f$  is not a convex function. ■

Theorem 22.5 can be strengthened to include nonopen sets by modifying the condition to be  $(\mathbf{y} - \mathbf{x})^\top \mathbf{F}(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq 0$  for all  $\mathbf{x}, \mathbf{y} \in \Omega$  (and assuming that  $f \in \mathcal{C}^2$  on some open set that contains  $\Omega$ ; for example,  $f \in \mathcal{C}^2$  on  $\mathbb{R}^n$ ). A proof similar to that above can be used in this case.

Note that by definition of concavity, a function  $f : \Omega \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^2$ , is concave over the convex set  $\Omega \subset \mathbb{R}^n$  if and only if for all  $\mathbf{x} \in \Omega$ , the Hessian  $\mathbf{F}(\mathbf{x})$  of  $f$  is negative semidefinite.

**Example 22.6** Determine whether the following functions are convex, concave, or neither:

1.  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = -8x^2$ .
2.  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f(\mathbf{x}) = 4x_1^2 + 3x_2^2 + 5x_3^2 + 6x_1x_2 + x_1x_3 - 3x_1 - 2x_2 + 15$ .
3.  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(\mathbf{x}) = 2x_1x_2 - x_1^2 - x_2^2$ .

*Solution:*

1. We use Theorem 22.5. We first compute the Hessian, which in this case is just the second derivative:  $(d^2f/dx^2)(x) = -16 < 0$  for all  $x \in \mathbb{R}$ . Hence,  $f$  is concave over  $\mathbb{R}$ .

2. The Hessian matrix of  $f$  is

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} 8 & 6 & 1 \\ 6 & 6 & 0 \\ 1 & 0 & 10 \end{bmatrix}.$$

The leading principal minors of  $\mathbf{F}(\mathbf{x})$  are

$$\Delta_1 = 8 > 0,$$

$$\Delta_2 = \det \begin{bmatrix} 8 & 6 \\ 6 & 6 \end{bmatrix} = 12 > 0,$$

$$\Delta_3 = \det \mathbf{F}(\mathbf{x}) = 114 > 0.$$

Hence,  $\mathbf{F}(\mathbf{x})$  is positive definite for all  $\mathbf{x} \in \mathbb{R}^3$ . Therefore,  $f$  is a convex function over  $\mathbb{R}^3$ .

3. The Hessian of  $f$  is

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} -2 & 2 \\ 2 & -2 \end{bmatrix},$$

which is negative semidefinite for all  $\mathbf{x} \in \mathbb{R}^2$ . Hence,  $f$  is concave on  $\mathbb{R}^2$ . ■

## 22.3 Convex Optimization Problems

In this section we consider optimization problems where the objective function is a convex function and the constraint set is a convex set. We refer to such problems as *convex optimization problems* or *convex programming problems*. Optimization problems that can be classified as convex programming problems include linear programs and optimization problems with quadratic objective function and linear constraints. Convex programming problems are interesting for several reasons. Specifically, as we shall see, local minimizers are global for such problems. Furthermore, first-order necessary conditions become sufficient conditions for minimization.

Our first theorem below states that in convex programming problems, local minimizers are also global.

**Theorem 22.6** *Let  $f : \Omega \rightarrow \mathbb{R}$  be a convex function defined on a convex set  $\Omega \subset \mathbb{R}^n$ . Then, a point is a global minimizer of  $f$  over  $\Omega$  if and only if it is a local minimizer of  $f$ .* □

*Proof.*  $\Rightarrow$ : This is obvious.

$\Leftarrow$ : We prove this by contraposition. Suppose that  $\mathbf{x}^*$  is not a global minimizer of  $f$  over  $\Omega$ . Then, for some  $\mathbf{y} \in \Omega$ , we have  $f(\mathbf{y}) < f(\mathbf{x}^*)$ . By assumption, the function  $f$  is convex, and hence for all  $\alpha \in (0, 1)$ ,

$$f(\alpha\mathbf{y} + (1 - \alpha)\mathbf{x}^*) \leq \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}^*).$$

Because  $f(\mathbf{y}) < f(\mathbf{x}^*)$ , we have

$$\alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}^*) = \alpha(f(\mathbf{y}) - f(\mathbf{x}^*)) + f(\mathbf{x}^*) < f(\mathbf{x}^*).$$

Thus, for all  $\alpha \in (0, 1)$ ,

$$f(\alpha\mathbf{y} + (1 - \alpha)\mathbf{x}^*) < f(\mathbf{x}^*).$$

Hence, there exist points that are arbitrarily close to  $\mathbf{x}^*$  and have lower objective function value. For example, the sequence  $\{\mathbf{y}_n\}$  of points given by

$$\mathbf{y}_n = \frac{1}{n}\mathbf{y} + \left(1 - \frac{1}{n}\right)\mathbf{x}^*$$

converges to  $\mathbf{x}^*$ , and  $f(\mathbf{y}_n) < f(\mathbf{x}^*)$ . Hence,  $\mathbf{x}^*$  is not a local minimizer, which completes the proof. ■

We now show that the set of global optimizers is convex. For this, we need the following lemma.

**Lemma 22.1** *Let  $g : \Omega \rightarrow \mathbb{R}$  be a convex function defined on a convex set  $\Omega \subset \mathbb{R}^n$ . Then, for each  $c \in \mathbb{R}$ , the set*

$$\Gamma_c = \{\mathbf{x} \in \Omega : g(\mathbf{x}) \leq c\}$$

*is a convex set.*

□

*Proof.* Let  $\mathbf{x}, \mathbf{y} \in \Gamma_c$ . Then,  $g(\mathbf{x}) \leq c$  and  $g(\mathbf{y}) \leq c$ . Because  $g$  is convex, for all  $\alpha \in (0, 1)$ ,

$$g(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha g(\mathbf{x}) + (1 - \alpha)g(\mathbf{y}) \leq c.$$

Hence,  $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \Gamma_c$ , which implies that  $\Gamma_c$  is convex. ■

**Corollary 22.1** *Let  $f : \Omega \rightarrow \mathbb{R}$  be a convex function defined on a convex set  $\Omega \subset \mathbb{R}^n$ . Then, the set of all global minimizers of  $f$  over  $\Omega$  is a convex set.*

□

*Proof.* The result follows immediately from Lemma 22.1 by setting

$$c = \min_{\mathbf{x} \in \Omega} f(\mathbf{x}).$$

■

We now show that if the objective function is continuously differentiable and convex, then the first-order necessary condition (see Theorem 6.1) for a point to be a minimizer is also sufficient. We use the following lemma.

**Lemma 22.2** *Let  $f : \Omega \rightarrow \mathbb{R}$  be a convex function defined on the convex set  $\Omega \subset \mathbb{R}^n$  and  $f \in \mathcal{C}^1$  on an open convex set containing  $\Omega$ . Suppose that the point  $\mathbf{x}^* \in \Omega$  is such that for all  $\mathbf{x} \in \Omega$ ,  $\mathbf{x} \neq \mathbf{x}^*$ , we have*

$$Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \geq 0.$$

*Then,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ .*  $\square$

*Proof.* Because the function  $f$  is convex, by Theorem 22.4, for all  $\mathbf{x} \in \Omega$ , we have

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

Hence, the condition  $Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \geq 0$  implies that  $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ .  $\blacksquare$

Observe that for any  $\mathbf{x} \in \Omega$ , the vector  $\mathbf{x} - \mathbf{x}^*$  can be interpreted as a feasible direction at  $\mathbf{x}^*$  (see Definition 6.2). Using Lemma 22.2, we have the following theorem (cf. Theorem 6.1).

**Theorem 22.7** *Let  $f : \Omega \rightarrow \mathbb{R}$  be a convex function defined on the convex set  $\Omega \subset \mathbb{R}^n$ , and  $f \in \mathcal{C}^1$  on an open convex set containing  $\Omega$ . Suppose that the point  $\mathbf{x}^* \in \Omega$  is such that for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have*

$$\mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0.$$

*Then,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ .*  $\square$

*Proof.* Let  $\mathbf{x} \in \Omega$ ,  $\mathbf{x} \neq \mathbf{x}^*$ . By convexity of  $\Omega$ ,

$$\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*) = \alpha\mathbf{x} + (1 - \alpha)\mathbf{x}^* \in \Omega$$

for all  $\alpha \in (0, 1)$ . Hence, the vector  $\mathbf{d} = \mathbf{x} - \mathbf{x}^*$  is a feasible direction at  $\mathbf{x}^*$  (see Definition 6.2). By assumption,

$$Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) = \mathbf{d}^\top \nabla f(\mathbf{x}^*) \geq 0.$$

Hence, by Lemma 22.2,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ .  $\blacksquare$

From Theorem 22.7, we easily deduce the following corollary (compare this with Corollary 6.1).

**Corollary 22.2** *Let  $f : \Omega \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , be a convex function defined on the convex set  $\Omega \subset \mathbb{R}^n$ . Suppose that the point  $\mathbf{x}^* \in \Omega$  is such that*

$$\nabla f(\mathbf{x}^*) = \mathbf{0}.$$

Then,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . □

We now consider the constrained optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}. \end{aligned}$$

We assume that the feasible set is convex. An example where this is the case is when

$$\mathbf{h}(\mathbf{x}) = A\mathbf{x} - \mathbf{b}.$$

The following theorem states that provided the feasible set is convex, the Lagrange condition is sufficient for a point to be a minimizer.

**Theorem 22.8** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , be a convex function on the set of feasible points*

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}\},$$

*where  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{h} \in \mathcal{C}^1$ , and  $\Omega$  is convex. Suppose that there exist  $\mathbf{x}^* \in \Omega$  and  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that*

$$Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) = \mathbf{0}^\top.$$

*Then,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ .* □

*Proof.* By Theorem 22.4, for all  $\mathbf{x} \in \Omega$ , we have

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

Substituting  $Df(\mathbf{x}^*) = -\boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*)$  into the inequality above yields

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) - \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

Because  $\Omega$  is convex,  $(1 - \alpha)\mathbf{x}^* + \alpha\mathbf{x} \in \Omega$  for all  $\alpha \in (0, 1)$ . Thus,

$$\mathbf{h}(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) = \mathbf{h}((1 - \alpha)\mathbf{x}^* + \alpha\mathbf{x}) = \mathbf{0}$$

for all  $\alpha \in (0, 1)$ . Premultiplying by  $\boldsymbol{\lambda}^{*\top}$ , subtracting  $\boldsymbol{\lambda}^{*\top} \mathbf{h}(\mathbf{x}^*) = 0$ , and dividing by  $\alpha$ , we get

$$\frac{\boldsymbol{\lambda}^{*\top} \mathbf{h}(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) - \boldsymbol{\lambda}^{*\top} \mathbf{h}(\mathbf{x}^*)}{\alpha} = 0$$

for all  $\alpha \in (0, 1)$ . If we now take the limit as  $\alpha \rightarrow 0$  and apply the definition of the directional derivative of  $\boldsymbol{\lambda}^{*\top} \mathbf{h}$  at  $\mathbf{x}^*$  in the direction  $\mathbf{x} - \mathbf{x}^*$  (see Section 6.2), we get

$$\boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) = 0.$$

Hence,

$$f(\mathbf{x}) \geq f(\mathbf{x}^*),$$

which implies that  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . ■

Consider the general constrained optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0} \\ & && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

As before, we assume that the feasible set is convex. This is the case if, for example, the two sets  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$  and  $\{\mathbf{x} : \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$  are convex, because the feasible set is the intersection of these two sets (see also Theorem 4.1). We have already seen an example where the set  $\{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$  is convex. On the other hand, an example where the set  $\{\mathbf{x} : \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$  is convex is when the components of  $\mathbf{g} = [g_1, \dots, g_p]^\top$  are all convex functions. Indeed, the set  $\{\mathbf{x} : \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\}$  is the intersection of the sets  $\{\mathbf{x} : g_i(\mathbf{x}) \leq 0\}$ ,  $i = 1, \dots, p$ . Because each of these sets is convex (by Lemma 22.1), their intersection is also convex.

We now prove that the Karush-Kuhn-Tucker (KKT) condition is sufficient for a point to be a minimizer to the problem above.

**Theorem 22.9** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , be a convex function on the set of feasible points*

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\},$$

where  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ ,  $\mathbf{h}, \mathbf{g} \in \mathcal{C}^1$ , and  $\Omega$  is convex. Suppose that there exist  $\mathbf{x}^* \in \Omega$ ,  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ , and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$ , such that

1.  $\boldsymbol{\mu}^* \geq \mathbf{0}$ .
2.  $Df(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top$ .
3.  $\boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = 0$ .

Then,  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ . □

*Proof.* Suppose that  $\mathbf{x} \in \Omega$ . By convexity of  $f$  and Theorem 22.4,

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + Df(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

Using condition 2, we get

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) - \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) - \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*).$$

As in the proof of Theorem 22.8, we can show that  $\boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) = 0$ . We now claim that  $\boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \leq 0$ . To see this, note that because  $\Omega$  is convex,  $(1 - \alpha)\mathbf{x}^* + \alpha\mathbf{x} \in \Omega$  for all  $\alpha \in (0, 1)$ . Thus,

$$\mathbf{g}(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) = \mathbf{g}((1 - \alpha)\mathbf{x}^* + \alpha\mathbf{x}) \leq \mathbf{0}$$

for all  $\alpha \in (0, 1)$ . Premultiplying by  $\mu^{*\top} \geq \mathbf{0}^\top$  (by condition 1), subtracting  $\mu^{*\top} g(\mathbf{x}^*) = 0$  (by condition 3), and dividing by  $\alpha$ , we get

$$\frac{\mu^{*\top} g(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) - \mu^{*\top} g(\mathbf{x}^*)}{\alpha} \leq 0.$$

We now take the limit as  $\alpha \rightarrow 0$  to obtain  $\mu^{*\top} Dg(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \leq 0$ .

From the above, we have

$$\begin{aligned} f(\mathbf{x}) &\geq f(\mathbf{x}^*) - \lambda^{*\top} Dh(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) - \mu^{*\top} Dg(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \\ &\geq f(\mathbf{x}^*) \end{aligned}$$

for all  $\mathbf{x} \in \Omega$ , which completes the proof. ■

**Example 22.7** A bank account starts out with 0 dollars. At the beginning of each month, we deposit some money into the bank account. Denote by  $x_k$  the amount deposited in the  $k$ th month,  $k = 1, 2, \dots$ . Suppose that the monthly interest rate is  $r > 0$  and the interest is paid into the account at the end of each month (and compounded). We wish to maximize the total amount of money accumulated at the end of  $n$  months, such that the total money deposited during the  $n$  months does not exceed  $D$  dollars (where  $D > 0$ ).

To solve this problem we first show that the problem can be posed as a linear program, and is therefore a convex optimization problem. Let  $y_k$  be the total amount in the bank at the end of the  $k$ th month. Then,  $y_k = (1+r)(y_{k-1} + x_k)$ ,  $k \geq 1$ , with  $y_0 = 0$ . Therefore, we want to maximize  $y_n$  subject to the constraint that  $x_k \geq 0$ ,  $k = 1, \dots, n$ , and  $x_1 + \dots + x_n \leq D$ . It is easy to deduce that

$$y_n = (1+r)^n x_1 + (1+r)^{n-1} x_2 + \dots + (1+r) x_n.$$

Let  $\mathbf{c}^\top = [(1+r)^n, (1+r)^{n-1}, \dots, (1+r)]$ ,  $\mathbf{e}^\top = [1, \dots, 1]$ , and  $\mathbf{x} = [x_1, \dots, x_n]^\top$ . Then, we can write the problem as

$$\begin{aligned} &\text{maximize } \mathbf{c}^\top \mathbf{x} \\ &\text{subject to } \mathbf{e}^\top \mathbf{x} \leq D \\ &\quad \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

which is a linear program.

It is intuitively clear that the optimal strategy is to deposit  $D$  dollars in the first month. To show that this strategy is indeed optimal, we use Theorem 22.9. Let  $\mathbf{x}^* = [D, 0, \dots, 0]^\top \in \mathbb{R}^n$ . Because the problem is a convex programming problem, it suffices to show that  $\mathbf{x}^*$  satisfies the KKT

condition (see Theorem 22.9). The KKT condition for this problem is

$$\begin{aligned} -\mathbf{c}^\top + \mu^{(1)} \mathbf{e}^\top - \boldsymbol{\mu}^{(2)\top} &= 0, \\ \mu^{(1)}(\mathbf{e}^\top \mathbf{x}^* - D) &= 0, \\ \boldsymbol{\mu}^{(2)\top} \mathbf{x}^* &= 0, \\ \mathbf{e}^\top \mathbf{x}^* - D &\leq 0, \\ -\mathbf{x}^* &\leq \mathbf{0}, \\ \mu^{(1)} &\geq 0, \\ \boldsymbol{\mu}^{(2)} &\geq \mathbf{0}, \\ \mathbf{e}^\top \mathbf{x} &\leq D, \\ \mathbf{x} &\geq \mathbf{0}, \end{aligned}$$

where  $\mu^{(1)} \in \mathbb{R}$  and  $\boldsymbol{\mu}^{(2)} \in \mathbb{R}^n$ . Let  $\mu^{(1)} = (1+r)^n$  and  $\boldsymbol{\mu}^{(2)} = (1+r)^n \mathbf{e} - \mathbf{c}$ . Then, it is clear that the KKT condition is satisfied. Therefore,  $\mathbf{x}^*$  is a global minimizer. ■

An entire book devoted to the vast topic of convexity and optimization is [7]. For extensions of the theory of convex optimization, we refer the reader to [136, Chapter 10]. The study of convex programming problems also serves as a prerequisite to *nondifferentiable optimization* (see, e.g., [38]).

## 22.4 Semidefinite Programming

Semidefinite programming is a subfield of convex optimization concerned with minimizing a linear objective function subject to a linear matrix inequality. The linear matrix inequality constraint defines a convex feasible set over which the linear objective function is to be minimized. Semidefinite programming can be viewed as an extension of linear programming, where the componentwise inequalities on vectors are replaced by matrix inequalities (see Exercise 22.20). For further reading on the subject of semidefinite programming, we recommend an excellent survey paper by Vandenberghe and Boyd [128].

### Linear Matrix Inequalities and Their Properties

Consider  $n+1$  real symmetric matrices

$$\mathbf{F}_i = \mathbf{F}_i^\top \in \mathbb{R}^{m \times m}, \quad i = 0, 1, \dots, n$$

and a vector

$$\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n.$$

Then,

$$\begin{aligned}\mathbf{F}(\mathbf{x}) &= \mathbf{F}_0 + x_1 \mathbf{F}_1 + \cdots + x_n \mathbf{F}_n \\ &= \mathbf{F}_0 + \sum_{i=1}^n x_i \mathbf{F}_i\end{aligned}$$

is an affine function of  $\mathbf{x}$ , because  $\mathbf{F}(\mathbf{x})$  is composed of a linear term  $\sum_{i=1}^n x_i \mathbf{F}_i$  and a constant term  $\mathbf{F}_0$ .

Consider now an inequality constraint of the form

$$\mathbf{F}(\mathbf{x}) = \mathbf{F}_0 + x_1 \mathbf{F}_1 + \cdots + x_n \mathbf{F}_n \geq 0.$$

The inequality constraint above is to be interpreted as the set of vectors  $\mathbf{x}$  such that

$$\mathbf{z}^\top \mathbf{F}(\mathbf{x}) \mathbf{z} \geq 0 \text{ for all } \mathbf{z} \in \mathbb{R}^m;$$

that is,  $\mathbf{F}(\mathbf{x})$  is positive semidefinite [or, in the usual notation,  $\mathbf{F}(\mathbf{x}) \geq 0$ ]. Recall that the terms  $\mathbf{F}_i$  represent constant matrices,  $\mathbf{x}$  is unknown, and  $\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{x})^\top$  is an affine function  $\mathbf{x}$ . The expression  $\mathbf{F}(\mathbf{x}) = \mathbf{F}_0 + x_1 \mathbf{F}_1 + \cdots + x_n \mathbf{F}_n \geq 0$  is referred to in the literature as a *linear matrix inequality* (LMI), although the term *affine matrix inequality* would seem to be more appropriate. It is easy to verify that the set  $\{\mathbf{x} : \mathbf{F}(\mathbf{x}) \geq 0\}$  is convex (see Exercise 22.20).

We can speak similarly of LMIs of the form  $\mathbf{F}(\mathbf{x}) > 0$ , where the requirement is for  $\mathbf{F}(\mathbf{x})$  to be positive definite (rather than just positive semidefinite). It is again easy to see that the set  $\{\mathbf{x} : \mathbf{F}(\mathbf{x}) > 0\}$  is convex.

A system of LMIs

$$\mathbf{F}_1(\mathbf{x}) \geq 0, \mathbf{F}_2(\mathbf{x}) \geq 0, \dots, \mathbf{F}_k(\mathbf{x}) \geq 0$$

can be represented as one single LMI:

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} \mathbf{F}_1(\mathbf{x}) & & & \\ & \mathbf{F}_2(\mathbf{x}) & & \\ & & \ddots & \\ & & & \mathbf{F}_k(\mathbf{x}) \end{bmatrix} \geq 0.$$

As an example, a linear inequality involving an  $m \times n$  real constant matrix  $\mathbf{A}$  of the form

$$\mathbf{A}\mathbf{x} \leq \mathbf{b}$$

can be represented as  $m$  LMIs:

$$b_i - \mathbf{a}_i^\top \mathbf{x} \geq 0, \quad i = 1, 2, \dots, m,$$

where  $\mathbf{a}_i^\top$  is the  $i$ th row of the matrix  $\mathbf{A}$ . We can view each scalar inequality as an LMI. We then represent  $m$  LMIs as one LMI:

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} b_1 - \mathbf{a}_1^\top \mathbf{x} \\ & b_2 - \mathbf{a}_2^\top \mathbf{x} \\ & & \ddots \\ & & & b_m - \mathbf{a}_m^\top \mathbf{x} \end{bmatrix} \geq 0.$$

With the foregoing facts as background, we can now give an example of semidefinite programming:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{F}(\mathbf{x}) \geq 0. \end{aligned}$$

The matrix property that we discuss next is useful when converting certain LMIs or nonlinear matrix inequalities into equivalent LMIs. We start with a simple observation. Let  $\mathbf{P}$  be a nonsingular  $n \times n$  matrix and let  $\mathbf{x} = \mathbf{M}\mathbf{z}$ , where  $\mathbf{M} \in \mathbb{R}^{n \times n}$  such that  $\det \mathbf{M} \neq 0$ . Then, we have

$$\mathbf{x}^\top \mathbf{P} \mathbf{x} \geq 0 \text{ if and only if } \mathbf{z}^\top \mathbf{M}^\top \mathbf{P} \mathbf{M} \mathbf{z} \geq 0;$$

that is,

$$\mathbf{P} \geq 0 \text{ if and only if } \mathbf{M}^\top \mathbf{P} \mathbf{M} \geq 0.$$

Similarly,

$$\mathbf{P} > 0 \text{ if and only if } \mathbf{M}^\top \mathbf{P} \mathbf{M} > 0.$$

Suppose that we have a square matrix

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix}.$$

Then, by the observation above,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix} \geq 0 \text{ if and only if } \begin{bmatrix} \mathbf{O} & \mathbf{I} \\ \mathbf{I} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{O} & \mathbf{I} \\ \mathbf{I} & \mathbf{O} \end{bmatrix} \geq 0,$$

where  $\mathbf{I}$  is an identity matrix of appropriate dimension. In other words,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{bmatrix} \geq 0 \text{ if and only if } \begin{bmatrix} \mathbf{D} & \mathbf{B}^\top \\ \mathbf{B} & \mathbf{A} \end{bmatrix} \geq 0.$$

We now introduce the notion of the Schur complement, useful in studying LMIs. Consider a square matrix of the form

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix},$$

where  $\mathbf{A}_{11}$  and  $\mathbf{A}_{22}$  are square submatrices. Suppose that the matrix  $\mathbf{A}_{11}$  is invertible. Then, we have

$$\begin{bmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{O} \\ \mathbf{O} & \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \end{bmatrix}.$$

Let

$$\Delta_{11} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12},$$

which is called the *Schur complement* of  $\mathbf{A}_{11}$ . For the case when  $\mathbf{A}_{12} = \mathbf{A}_{21}^\top$ , we have

$$\begin{bmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{21}^\top \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{A}_{11}^{-1}\mathbf{A}_{21}^\top \\ \mathbf{O} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{O} \\ \mathbf{O} & \Delta_{11} \end{bmatrix},$$

where

$$\Delta_{11} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{21}^\top.$$

Hence,

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{21}^\top \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} > 0 \text{ if and only if } \begin{bmatrix} \mathbf{A}_{11} & \mathbf{O} \\ \mathbf{O} & \Delta_{11} \end{bmatrix} > 0;$$

that is,

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{21}^\top \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} > 0 \text{ if and only if } \mathbf{A}_{11} > 0 \text{ and } \Delta_{11} > 0.$$

Given

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix},$$

we can similarly define the Schur complement of  $\mathbf{A}_{22}$ , assuming that  $\mathbf{A}_{22}$  is invertible. We have

$$\begin{bmatrix} \mathbf{I} & -\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ -\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \Delta_{22} & \mathbf{O} \\ \mathbf{O} & \mathbf{A}_{22} \end{bmatrix},$$

where  $\Delta_{22} = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$  is the Schur complement of  $\mathbf{A}_{22}$ . So, for the case where  $\mathbf{A}_{12} = \mathbf{A}_{21}^\top$ ,

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{21}^\top \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} > 0 \text{ if and only if } \mathbf{A}_{22} > 0 \text{ and } \Delta_{22} > 0.$$

Many problems of optimization, control design, and signal processing can be formulated in terms of LMIs. To determine whether or not there exists a point  $\mathbf{x}$  such that  $\mathbf{F}(\mathbf{x}) > 0$  is called a *feasibility problem*. We say that the LMI is nonfeasible if no such solution exists.

**Example 22.8** We now present a simple example illustrating the LMI feasibility problem. Let  $\mathbf{A} \in \mathbb{R}^{m \times m}$  be a given real constant square matrix. Suppose that we wish to determine if  $\mathbf{A}$  has all its eigenvalues in the open left half-complex plane. It is well known that this condition is true if and only if there exists a real symmetric positive definite matrix  $\mathbf{P}$  such that

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0,$$

or, equivalently,

$$-\mathbf{A}^\top \mathbf{P} - \mathbf{P} \mathbf{A} > 0$$

(also called the *Lyapunov inequality*; see [16]). Thus, the location of all eigenvalues of  $\mathbf{A}$  being in the open left half-complex plane is equivalent to feasibility of the following matrix inequality:

$$\begin{bmatrix} \mathbf{P} & \mathbf{O} \\ \mathbf{O} & -\mathbf{A}^\top \mathbf{P} - \mathbf{P} \mathbf{A} \end{bmatrix} > 0;$$

that is, the existence of  $\mathbf{P} = \mathbf{P}^\top > 0$  such that  $\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0$ .

We now show that finding  $\mathbf{P} = \mathbf{P}^\top > 0$  such that  $\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0$  is indeed an LMI. For this, let

$$\mathbf{P} = \begin{bmatrix} x_1 & x_2 & \cdots & x_m \\ x_2 & x_{m+1} & \cdots & x_{2m-1} \\ \vdots & & & \vdots \\ x_m & x_{2m-1} & \cdots & x_n \end{bmatrix},$$

where

$$n = \frac{m(m+1)}{2}.$$

We next define the following matrices:

$$\begin{aligned}\mathbf{P}_1 &= \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \\ \mathbf{P}_2 &= \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \\ &\vdots \\ \mathbf{P}_n &= \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix}.\end{aligned}$$

Note that  $\mathbf{P}_i$  has only nonzero elements corresponding to  $\mathbf{x}_i$  in  $\mathbf{P}$ . Let

$$\mathbf{F}_i = -\mathbf{A}^\top \mathbf{P}_i - \mathbf{P}_i \mathbf{A}, \quad i = 1, 2, \dots, n.$$

We can then write

$$\begin{aligned}\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} &= x_1 (\mathbf{A}^\top \mathbf{P}_1 + \mathbf{P}_1 \mathbf{A}) + x_2 (\mathbf{A}^\top \mathbf{P}_2 + \mathbf{P}_2 \mathbf{A}) + \cdots \\ &\quad + x_n (\mathbf{A}^\top \mathbf{P}_n + \mathbf{P}_n \mathbf{A}) \\ &= -x_1 \mathbf{F}_1 - x_2 \mathbf{F}_2 - \cdots - x_n \mathbf{F}_n \\ &< 0.\end{aligned}$$

Let

$$\mathbf{F}(\mathbf{x}) = x_1 \mathbf{F}_1 + x_2 \mathbf{F}_2 + \cdots + x_n \mathbf{F}_n.$$

Then,

$$\mathbf{P} = \mathbf{P}^\top > 0 \quad \text{and} \quad \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0$$

if and only if

$$\mathbf{F}(\mathbf{x}) > 0.$$

Note that this LMI involves a strict inequality. Most numerical solvers do not handle strict inequalities. Such solvers simply treat a strict inequality ( $>$ ) as a non-strict inequality ( $\geq$ ). ■

## LMI Solvers

The inequality  $\mathbf{F}(\mathbf{x}) = \mathbf{F}_0 + x_1\mathbf{F}_1 + \cdots + x_n\mathbf{F}_n \geq 0$  is called the *canonical representation* of an LMI. Numerical LMI solvers do not deal directly with LMIs in canonical form because of various inefficiencies. Instead, LMI solvers use a structured representation of LMIs.

We can use MATLAB's LMI toolbox to solve LMIs efficiently. This toolbox has three types of LMI solvers, which we discuss next.

### Finding a Feasible Solution Under LMI Constraints

First, we discuss MATLAB's LMI solver for solving the feasibility problem defined by a given system of LMI constraints. Using this solver, we can solve any system of LMIs of the form

$$\mathbf{N}^\top \mathcal{L}(\mathbf{X}_1, \dots, \mathbf{X}_k) \mathbf{N} \leq \mathbf{M}^\top \mathcal{R}(\mathbf{X}_1, \dots, \mathbf{X}_k) \mathbf{M},$$

where  $\mathbf{X}_1, \dots, \mathbf{X}_k$  are matrix variables,  $\mathbf{N}$  is the left outer factor,  $\mathbf{M}$  is the right outer factor,  $\mathcal{L}(\mathbf{X}_1, \dots, \mathbf{X}_k)$  is the left inner factor, and  $\mathcal{R}(\mathbf{X}_1, \dots, \mathbf{X}_k)$  is the right inner factor. The matrices  $\mathcal{L}(\cdot)$  and  $\mathcal{R}(\cdot)$  are, in general, symmetric block matrices. We note that the term left-hand side refers to what is on the “smaller” side of the inequality  $0 \leq \mathbf{X}$ . Thus in  $\mathbf{X} \geq 0$ , the matrix  $\mathbf{X}$  is still on the right-hand side because it is on the “larger” side of the inequality.

We now provide a description of an approach that can be used to solve the given LMI system feasibility problem. To initialize the LMI system description, we type `setlmis([])`. Then we declare matrix variables using the command `lmivar`. The command `lmitemr` allows us to specify LMIs that constitute the LMI system under consideration. Next, we need to obtain an internal representation using the command `getlmis`. We next compute a feasible solution to the LMI system using the command `feasp`. After that, we extract matrix variable values with the command `dec2mat`. In summary, a general structure of a MATLAB program for finding a feasible solution to the set of LMIs could have the form

```
setlmis([])
lmivar
lmitemr
.
.
.
lmitemr
getlmis
feasp
dec2mat
```

We now analyze these commands in some detail so that the reader can write simple MATLAB programs for solving LMIs after completing this section.

First, to create a new matrix-valued variable, say,  $\mathbf{X}$ , in the given LMI system, we use the command

```
X = lmivar(type,structure)
```

The input **type** specifies the structure of the variable  $\mathbf{X}$ . There may be three structures of matrix variables. When **type**=1, we have a symmetric block diagonal matrix variable. The input **type**=2 refers to a full rectangular matrix variable. Finally, **type**=3 refers to other cases. The second input **structure** gives additional information on the structure of the matrix variable  $\mathbf{X}$ . For example, the matrix variable  $\mathbf{X}$  could have the form

$$\mathbf{X} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{D}_2 & \cdots & \mathbf{O} \\ \vdots & & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{D}_r \end{bmatrix},$$

where each  $\mathbf{D}_i$  is a square symmetric matrix. For the example above we would use **type**=1. The matrix variable above has  $r$  blocks. The input **structure** is then an  $r \times 2$  matrix whose  $i$ th row describes the  $i$ th block, where the first component of each row gives the corresponding block size, while the second element of each row specifies the block type. For example,

```
X = lmivar(1,[3 1])
```

specifies a full symmetric  $3 \times 3$  matrix variable. On the other hand,

```
X = lmivar(2,[2 3])
```

specifies a rectangular  $2 \times 3$  matrix variable. Finally, a matrix variable  $\mathbf{S}$  of the form

$$\mathbf{S} = \begin{bmatrix} s_1 & 0 & | & 0 & 0 \\ 0 & s_1 & | & 0 & 0 \\ --- & --- & | & --- & --- \\ 0 & 0 & | & s_2 & s_3 \\ 0 & 0 & | & s_3 & s_4 \end{bmatrix}$$

can be declared as follows:

```
S = lmivar(1,[2 0;2 1])
```

Note above that the second component of the first row of the second input has the value of zero; that is, **structure**(1,2)=0. This describes a scalar block matrix of the form

$$\mathbf{D}_1 = s_1 \mathbf{I}_2.$$

Note that the second block is a  $2 \times 2$  symmetric full block.

We next take a closer look at a command whose purpose is to specify the terms of the LMI system of interest. This command has the form

```
lmitem(termid,A,B,flag)
```

We briefly describe each of the four inputs of this command. The first input, `termid`, is a row with four elements that specify the terms of each LMI of the LMI system. We have `termid(1)=n` to specify the left-hand side of the  $n$ th LMI. We use `termid(1)=-n` to specify the right-hand side of the  $n$ th LMI. The middle two elements of the input `termid` specify the block location. Thus `termid(2,3)=[i j]` refers to the term that belongs to the  $(i, j)$  block of the LMI specified by the first component. Finally, `termid(4)=0` for the constant term, `termid(4)=X` for the variable term in the form  $\mathbf{A} \mathbf{X} \mathbf{B}$ , while `termid(4)=-X` for the variable term in the form  $\mathbf{A} \mathbf{X}^\top \mathbf{B}$ . The second and third inputs of the command `lmitem` give the values of the left and right outer factors; that is, `A` and `B` give the values of the constant outer factors in the variable terms  $\mathbf{A} \mathbf{X} \mathbf{B}$  and  $\mathbf{A} \mathbf{X}^\top \mathbf{B}$ . Finally, the fourth input to `lmitem` serves as a compact way to specify the expression

$$\mathbf{A} \mathbf{X} \mathbf{B} + (\mathbf{A} \mathbf{X} \mathbf{B})^\top.$$

Thus, `flag='s'` can be used to denote a symmetrized expression. We now illustrate the command above on the following LMI:

$$\mathbf{P} \mathbf{A} + (\mathbf{P} \mathbf{A})^\top \leq 0.$$

We have one LMI with two terms. We could use the following description of this single LMI:

```
lmitem([1 1 1 P],1,A)
lmitem([1 1 1 -P],A',1)
```

On the other hand, we can describe this LMI compactly using the `flag` as follows:

```
lmitem([1 1 1 P],1,A,'s')
```

Now, to solve the feasibility problem we could have typed

```
[tmin,xfeas] = feas(lmis)
```

In general, for a given LMI feasibility problem of the form

$$\begin{aligned} &\text{find } \mathbf{x} \\ &\text{such that } \mathbf{L}(\mathbf{x}) \leq \mathbf{R}(\mathbf{x}), \end{aligned}$$

the command `feasp` solves the auxiliary convex problem

$$\begin{aligned} &\text{minimize } t \\ &\text{subject to } \mathbf{L}(\mathbf{x}) \leq \mathbf{R}(\mathbf{x}) + t\mathbf{I}. \end{aligned}$$

The system of LMIs is feasible if the minimal  $t$  is negative. We add that the current value of  $t$  is displayed by **feasp** at each iteration.

Finally, we convert the output of the LMI solver into matrix variables using the command

```
P = dec2mat(lmis,xfeas,P).
```

**Example 22.9** Let

$$\mathbf{A}_1 = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} \quad \text{and} \quad \mathbf{A}_2 = \begin{bmatrix} -2 & 0 \\ 1 & -1 \end{bmatrix}.$$

We use the commands of the LMI Control Toolbox discussed above to write a program that finds  $\mathbf{P}$  such that  $\mathbf{P} \geq 0.5\mathbf{I}_2$  and

$$\begin{aligned}\mathbf{A}_1^\top \mathbf{P} + \mathbf{P} \mathbf{A}_1 &\leq 0, \\ \mathbf{A}_2^\top \mathbf{P} + \mathbf{P} \mathbf{A}_2 &\leq 0.\end{aligned}$$

The program is as follows:

```
A_1 = [-1 0;0 -1];
A_2 = [-2 0;1 -1];
setlmis([])
P = lmivar(1,[2,1])
lmiterm([1 1 1 P],A_1',1,'s')
lmiterm([2 1 1 P],A_2',1,'s')
lmiterm([3 1 1 0],.5)
lmiterm([-3 1 1 P],1,1)
lmis=getlmis;
[tmin,xfeas] = feasp(lmis);
P = dec2mat(lmis,xfeas,P)
```

■

### Minimizing a Linear Objective Under LMI Constraints

The next solver we discuss solves the convex optimization problem

$$\begin{aligned}&\text{minimize} \quad \mathbf{c}^\top \mathbf{x} \\ &\text{subject to} \quad \mathbf{A}(\mathbf{x}) \leq \mathbf{B}(\mathbf{x}).\end{aligned}$$

The notation  $\mathbf{A}(\mathbf{x}) \leq \mathbf{B}(\mathbf{x})$  is shorthand notation for a general structured LMI system.

This solver is invoked using the function **mincx**. Thus, to solve a **mincx** problem, in addition to specifying the LMI constraints as in the **feasp** problem, we also declare the linear objective function. Then we invoke the function

`mincx`. We illustrate and contrast the `feasp` and `mincx` solvers in the following example.

**Example 22.10** Consider the optimization problem

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{c}^\top &= \begin{bmatrix} 4 & 5 \end{bmatrix}, \\ \mathbf{A} &= \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 2 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 8 \\ 18 \\ 14 \end{bmatrix}. \end{aligned}$$

We first solve the feasibility problem; that is, we find an  $\mathbf{x}$  such that  $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ , using the `feasp` solver. After that, we solve the minimization problem above using the `mincx` solver. A simple MATLAB code accomplishing these tasks is shown below.

```
% Enter problem data
A = [1 1; 1 3; 2 1];
b = [8 18 14]';
c = [-4 -5]';
setlmis([]);
X = lmivar(2,[2 1]);
lmiterm([1 1 1 X],A(1,:),1);
lmiterm([1 1 1 0],-b(1));
lmiterm([1 2 2 X],A(2,:),1);
lmiterm([1 2 2 0],-b(2));
lmiterm([1 3 3 X],A(3,:),1);
lmiterm([1 3 3 0],-b(3));
lmis = getlmis;
%-----
disp('-----feasp result-----')
[tmin,xfeas] = feasp(lmis);
x_feasp = dec2mat(lmis,xfeas,X)
disp('-----mincx result-----')
[objective,x_mincx] = mincx(lmis,c,[0.0001 1000 0 0 1])
```

The `feasp` function produces

$$\mathbf{x}_{\text{feasp}} = \begin{bmatrix} -64.3996 \\ -25.1712 \end{bmatrix}.$$

The `mincx` function produces

$$\mathbf{x}_{\text{mincx}} = \begin{bmatrix} 3.0000 \\ 5.0000 \end{bmatrix}.$$

■

In the next example, we discuss the function `defcx`, which we can use to construct the vector  $\mathbf{c}$  used by the LMI solver `mincx`.

**Example 22.11** Suppose that we wish to solve the optimization problem

$$\begin{aligned} & \text{minimize} && \text{trace}(\mathbf{P}) \\ & \text{subject to} && \mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} \leq 0 \\ & && \mathbf{P} \geq 0 \end{aligned}$$

where  $\text{trace}(\mathbf{P})$  is the sum of the diagonal elements of  $\mathbf{P}$ . We can use the function `mincx` to solve this problem. However, to use `mincx`, we need a vector  $\mathbf{c}$  such that

$$\mathbf{c}^\top \mathbf{x} = \text{trace}(\mathbf{P}).$$

After specifying the LMIs and obtaining their internal representation using, for example, the command `lmisys=getlmis`, we can obtain the desired  $\mathbf{c}$  with the following MATLAB code,

```
q = decnbr(lmisys);
c = zeros(q,1);
for j = 1:q
    Pj = defcx(lmisys,j,P);
    c(j) = trace(Pj);
end
```

Having obtained the vector  $\mathbf{c}$ , we can use the function `mincx` to solve the optimization problem. ■

### Minimizing a Generalized Eigenvalue Under LMI Constraints

This problem can be stated as

$$\begin{aligned} & \text{minimize} && \lambda \\ & \text{subject to} && \mathbf{C}(\mathbf{x}) \leq \mathbf{D}(\mathbf{x}) \\ & && 0 \leq \mathbf{B}(\mathbf{x}) \\ & && \mathbf{A}(\mathbf{x}) \leq \lambda \mathbf{B}(\mathbf{x}). \end{aligned}$$

Here, we need to distinguish between standard LMI constraints of the form  $\mathbf{C}(\mathbf{x}) \leq \mathbf{D}(\mathbf{x})$  and *linear-fractional LMIs* of the form  $\mathbf{A}(\mathbf{x}) \leq \lambda \mathbf{B}(\mathbf{x})$ , which

are concerned with the generalized eigenvalue  $\lambda$ . The generalized eigenvalue minimization problem under LMI constraints can be solved using the solver `gevp`. The basic structure of the `gevp` solver has the form

```
[lopt,xopt] = gevp{lmisys,nflc}
```

which returns `lopt`, the global minimum of the generalized eigenvalue, and `xopt`, the optimal decision vector variable. The argument `lmisys` is the system of LMIs,  $C(\mathbf{x}) \leq D(\mathbf{x})$ ,  $C(\mathbf{x}) \leq D(\mathbf{x})$ , and  $A(\mathbf{x}) \leq \lambda B(\mathbf{x})$  for  $\lambda = 1$ . As in the previous solvers, the corresponding optimal values of the matrix variables are obtained using `dec2mat`. The number of linear-fractional constraints is specified with `nflc`. There are other inputs to `gevp` but they are optional. For more information on this type of the LMI solver, we refer the reader to the LMI Lab in MATLAB's Robust Control Toolbox user's guide.

**Example 22.12** Consider the problem of finding the smallest  $\alpha$  such that

$$\begin{aligned} P &> 0 \\ A^T P + PA &\leq -\alpha P, \end{aligned}$$

where

$$A = \begin{bmatrix} -1.1853 & 0.9134 & 0.2785 \\ 0.9058 & -1.3676 & 0.5469 \\ 0.1270 & 0.0975 & -3.0000 \end{bmatrix}.$$

This problem is related to finding the decay rate of the stable linear differential equation  $\dot{\mathbf{x}} = A\mathbf{x}$ . Finding  $\alpha$  that solves the optimization problem above can be accomplished using the following LMIs:

```
A = [-1.1853 0.9134 0.2785
      0.9058 -1.3676 0.5469
      0.1270 0.0975 -3.0000];
setlmis([]);
P = lmivar(1,[3 1])
lmiterm([-1 1 1 P],1,1)    % P
lmiterm([1 1 1 0],.01)      % P >= 0.01*I
lmiterm([2 1 1 P],1,A,'s') % linear fractional constraint---LHS
lmiterm([-2 1 1 P],1,1)    % linear fractional constraint---RHS
lmis = getlmis;
[gamma,P_opt] = gevp(lmis,1);
P = dec2mat(lmis,P_opt,P)
alpha = -gamma
```

The result is

$$\alpha = 0.6561 \quad \text{and} \quad P = \begin{bmatrix} 0.6996 & -0.7466 & -0.0296 \\ -0.7466 & 0.8537 & -0.2488 \\ -0.0296 & -0.2488 & 3.2307 \end{bmatrix}.$$

Notice that we used  $\mathbf{P} \geq 0.01\mathbf{I}$  in place of  $\mathbf{P} > 0$ . ■

More examples of linear matrix inequalities in system and control theory can be found in the book by Boyd et al. [16].

A quick introduction to MATLAB's LMI toolbox is the tutorial that can be accessed with the command `lmdem` within MATLAB. In addition to the MATLAB's LMI toolbox, there is another toolbox for solving LMIs called LMITOOL, a built-in software package in Scilab toolbox, developed at INRIA in France. Scilab offers free software for numerical optimization. There is a version of LMITOOL for MATLAB that can be obtained from the website of the Scilab Consortium.

Yet Another LMI Package, YALMIP, for solving LMIs was developed in Switzerland in the Automatic Control Laboratory at ETH. YALMIP is an “intuitive and flexible modelling language for solving optimization problems in MATLAB.”

LMIs are tools of modern optimization. The following quote on numerical linear algebra from Gill, Murray, and Wright [52, p. 2] applies as well to the contents of this chapter: “At the heart of modern optimization methods are techniques associated with *linear algebra*. Numerical linear algebra applies not simply in optimization, but in all fields of scientific computation, including approximation, ordinary differential equations, and partial differential equations. *The importance of numerical linear algebra to modern scientific computing cannot be overstated*. Without fast and reliable linear algebraic building blocks, it is impossible to develop effective optimization methods; without some knowledge of the fundamental issues in linear algebra, it is impossible to understand what happens during the transition from equations in a textbook to actual computation.”

## EXERCISES

**22.1** Find the range of values of the parameter  $\alpha$  for which the function

$$f(x_1, x_2, x_3) = 2x_1x_3 - x_1^2 - x_2^2 - 5x_3^2 - 2\alpha x_1x_2 - 4x_2x_3$$

is concave.

**22.2** Consider the function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$  and  $\mathbf{x}, \mathbf{b} \in \mathbb{R}^n$ . Define the function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  by  $\phi(\alpha) = f(\mathbf{x} + \alpha\mathbf{d})$ , where  $\mathbf{x}, \mathbf{d} \in \mathbb{R}^n$  are fixed vectors and  $\mathbf{d} \neq \mathbf{0}$ . Show that  $\phi(\alpha)$  is a strictly convex quadratic function of  $\alpha$ .

**22.3** Show that  $f(\mathbf{x}) = x_1 x_2$  is a convex function on  $\Omega = \{[a, ma]^\top : a \in \mathbb{R}\}$ , where  $m$  is any given nonnegative constant.

**22.4** Suppose that the set  $\Omega = \{\mathbf{x} : h(\mathbf{x}) = c\}$  is convex, where  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $c \in \mathbb{R}$ . Show that  $h$  is convex and concave over  $\Omega$ .

**22.5** Find all subgradients of

$$f(x) = |x|, \quad x \in \mathbb{R},$$

at  $x = 0$  and at  $x = 1$ .

**22.6** Let  $\Omega \subset \mathbb{R}^n$  be a convex set, and  $f_i : \Omega \rightarrow \mathbb{R}$ ,  $i = 1, \dots, \ell$  be convex functions. Show that  $\max\{f_1, \dots, f_\ell\}$  is a convex function.

*Note:* The notation  $\max\{f_1, \dots, f_\ell\}$  denotes a function from  $\Omega$  to  $\mathbb{R}$  such that for each  $\mathbf{x} \in \Omega$ , its value is the largest among the numbers  $f_i(\mathbf{x})$ ,  $i = 1, \dots, \ell$ .

**22.7** Let  $\Omega \subset \mathbb{R}^n$  be an open convex set. Show that a symmetric matrix  $\mathbf{Q} \in \mathbb{R}^n$  is positive semidefinite if and only if for each  $\mathbf{x}, \mathbf{y} \in \Omega$ ,  $(\mathbf{x} - \mathbf{y})^\top \mathbf{Q} (\mathbf{x} - \mathbf{y}) \geq 0$ . Show that a similar result for positive definiteness holds if we replace the “ $\geq$ ” by “ $>$ ” in the inequality above.

**22.8** Consider the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \\ & \text{subject to} && x_1 + \cdots + x_n = 1 \\ & && x_1, \dots, x_n \geq 0 \end{aligned}$$

(see also Exercise 21.9). Is the problem a convex optimization problem? If yes, give a complete proof. If no, explain why not, giving complete explanations.

**22.9** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $f(\mathbf{x}) = x_1 x_2^2$ , where  $\mathbf{x} = [x_1, x_2]^\top$ , and  $\Omega = \{\mathbf{x} \in \mathbb{R}^2 : x_1 = x_2, x_1 \geq 0\}$ . (See also Exercise 21.8.) Show that the problem is a convex optimization problem.

**22.10** Consider the convex optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

Suppose that the points  $\mathbf{y} \in \Omega$  and  $\mathbf{z} \in \Omega$  are local minimizers. Determine the largest set of points  $G \subset \Omega$  for which you can be sure that every point in  $G$  is a global minimizer.

**22.11** Suppose that we have a convex optimization problem on  $\mathbb{R}^3$ .

- Consider the following three feasible points:  $[1, 0, 0]^\top, [0, 1, 0]^\top, [0, 0, 1]^\top$ . Suppose that all three have objective function value 1. What can you say about the objective function value of the point  $(1/3)[1, 1, 1]^\top$ ? Explain fully.
- Suppose we know that the three points in part a are global minimizers. What can you say about the point  $(1/3)[1, 1, 1]^\top$ ? Explain fully.

**22.12** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \\ & \text{subject to} && \mathbf{A} \mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{Q} = \mathbf{Q}^\top > 0$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and  $\text{rank } \mathbf{A} = m$ .

- Find all points satisfying the Lagrange condition for the problem (in terms of  $\mathbf{Q}$ ,  $\mathbf{A}$ , and  $\mathbf{b}$ ).
- Are the points (or point) global minimizers for this problem?

**22.13** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f \in \mathcal{C}^1$ , be a convex function on the set of feasible points

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{a}_i^\top \mathbf{x} + b_i \geq 0, \quad i = 1, \dots, p\},$$

where  $\mathbf{a}_1, \dots, \mathbf{a}_p \in \mathbb{R}^n$ , and  $b_1, \dots, b_p \in \mathbb{R}$ . Suppose that there exist  $\mathbf{x}^* \in S$ , and  $\boldsymbol{\mu}^* \in \mathbb{R}^p$ ,  $\boldsymbol{\mu}^* \leq \mathbf{0}$ , such that

$$Df(\mathbf{x}^*) + \sum_{j \in J(\mathbf{x}^*)} \mu_j^* \mathbf{a}_j^\top = \mathbf{0}^\top,$$

where  $J(\mathbf{x}^*) = \{i : \mathbf{a}_i^\top \mathbf{x}^* + b_i = 0\}$ . Show that  $\mathbf{x}^*$  is a global minimizer of  $f$  over  $\Omega$ .

**22.14** Consider the problem: minimize  $\|\mathbf{x}\|^2$  ( $\mathbf{x} \in \mathbb{R}^n$ ) subject to  $\mathbf{a}^\top \mathbf{x} \geq b$ , where  $\mathbf{a} \in \mathbb{R}^n$  is a nonzero vector and  $b \in \mathbb{R}$ ,  $b > 0$ . Suppose that  $\mathbf{x}^*$  is a solution to the problem.

- Show that the constraint set is convex.
- Use the KKT theorem to show that  $\mathbf{a}^\top \mathbf{x}^* = b$ .

- c. Show that  $\mathbf{x}^*$  is unique, and find an expression for  $\mathbf{x}^*$  in terms of  $\mathbf{a}$  and  $b$ .

**22.15** Consider the problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n \\ & \text{subject to} && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

For this problem we have the following theorem (see also Exercise 17.16).

**Theorem:** *A solution to this problem exists if and only if  $\mathbf{c} \geq \mathbf{0}$ . Moreover, if a solution exists,  $\mathbf{0}$  is a solution.*

- a. Show that the problem is a convex programming problem.
- b. Use the first-order necessary condition (for set constraints) to prove the theorem.
- c. Use the KKT condition to prove the above theorem.

**22.16** Consider a linear programming problem in standard form.

- a. Derive the KKT condition for the problem.
- b. Explain precisely why the KKT condition is sufficient for optimality in this case.
- c. Write down the dual to the standard form primal problem (see Chapter 17).
- d. Suppose that  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  are feasible solutions to the primal and dual, respectively. Use the KKT condition to prove that if the complementary slackness condition  $(\mathbf{c}^\top - \boldsymbol{\lambda}^{*\top} \mathbf{A})\mathbf{x}^* = 0$  holds, then  $\mathbf{x}^*$  is an optimal solution to the primal problem. Compare the above with Exercise 21.15.

**22.17** Consider two real-valued discrete-time signals,  $\mathbf{s}^{(1)}$  and  $\mathbf{s}^{(2)}$ , defined over the time interval  $[1, n]$ . Let  $s_i^{(1)}$  and  $s_i^{(2)}$  be the values at time  $i$  of the signals  $\mathbf{s}^{(1)}$  and  $\mathbf{s}^{(2)}$ , respectively. Assume that the energies of the two signals are 1 [i.e.,  $(s_1^{(1)})^2 + \dots + (s_n^{(1)})^2 = 1$  and  $(s_1^{(2)})^2 + \dots + (s_n^{(2)})^2 = 1$ ].

Let  $S_a$  be the set of all signals that are linear combinations of  $\mathbf{s}^{(1)}$  and  $\mathbf{s}^{(2)}$  with the property that for each signal in  $S_a$ , the value of the signal over all time is no smaller than  $a \in \mathbb{R}$ . For each  $\mathbf{s} \in S_a$ , if  $\mathbf{s} = x_1 \mathbf{s}^{(1)} + x_2 \mathbf{s}^{(2)}$ , we call  $x_1$  and  $x_2$  the *coefficients* of  $\mathbf{s}$ .

We wish to find a signal in  $S_a$  such that the sum of the squares of its coefficients is minimized.

- a. Formulate the problem as an optimization problem.

- b. Derive the Karush-Kuhn-Tucker conditions for the problem.
- c. Suppose that you have found a point satisfying the Karush-Kuhn-Tucker conditions. Does this point satisfy the second-order sufficient condition?
- d. Is this problem a convex optimization problem?

**22.18** Let a *probability vector* be any vector  $\mathbf{p} \in \mathbb{R}^n$  satisfying  $p_i > 0$ ,  $i = 1, \dots, n$ , and  $p_1 + \dots + p_n = 1$ .

Let  $\mathbf{p} \in \mathbb{R}^n$  and  $\mathbf{q} \in \mathbb{R}^n$  be two probability vectors. Define

$$D(\mathbf{p}, \mathbf{q}) = p_1 \log\left(\frac{p_1}{q_1}\right) + \dots + p_n \log\left(\frac{p_n}{q_n}\right),$$

where “log” is the natural logarithm function.

- a. Let  $\Omega$  be the set of all probability vectors (with fixed  $n$ ). Show that  $\Omega$  is convex.
- b. Show that for each fixed  $\mathbf{p}$ , the function  $f$  defined by  $f(\mathbf{q}) = D(\mathbf{p}, \mathbf{q})$  is convex over  $\Omega$ .
- c. Show the following:  $D(\mathbf{p}, \mathbf{q}) \geq 0$  for any probability vectors  $\mathbf{p}$  and  $\mathbf{q}$ . Moreover,  $D(\mathbf{p}, \mathbf{q}) = 0$  if and only if  $\mathbf{p} = \mathbf{q}$ .
- d. Describe an application of the result of part c.

**22.19** Let  $\Omega \subset \mathbb{R}^n$  be a nonempty closed convex set and  $\mathbf{z} \in \mathbb{R}^n$  be a given point such that  $\mathbf{z} \notin \Omega$ . Consider the optimization problem

$$\begin{aligned} &\text{minimize} && \|\mathbf{x} - \mathbf{z}\| \\ &\text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

Does this problem have an optimal solution? If so, is it unique? Whatever your assertion, prove it.

*Hint:* (i) If  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are optimal solutions, what can you say about  $\mathbf{x}_3 = (\mathbf{x}_1 + \mathbf{x}_2)/2$ ? (ii) The triangle inequality states that  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ , with equality holding if and only if  $\mathbf{x} = \alpha\mathbf{y}$  for some  $\alpha \geq 0$  (or  $\mathbf{x} = 0$  or  $\mathbf{y} = 0$ ).

**22.20** This exercise is about *semidefinite programming*.

- a. Show that if  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{B} \in \mathbb{R}^{n \times n}$  are symmetric and  $\mathbf{A} \geq 0$ ,  $\mathbf{B} \geq 0$ , then for any  $\alpha \in (0, 1)$ , we have  $\alpha\mathbf{A} + (1 - \alpha)\mathbf{B} \geq 0$ . As usual, the notation “ $\geq 0$ ” denotes positive semidefiniteness.

- b.** Consider the following semidefinite programming problem, that is, an optimization problem with linear objective function and linear matrix inequality constraints:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{F}_0 + \sum_{j=1}^n x_j \mathbf{F}_j \geq 0, \end{aligned}$$

where  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$  is the decision variable,  $\mathbf{c} \in \mathbb{R}^n$ , and  $\mathbf{F}_0, \mathbf{F}_1, \dots, \mathbf{F}_n \in \mathbb{R}^{m \times m}$  are symmetric.

Show that this problem is a convex optimization problem.

- c.** Consider the linear programming problem

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and the inequality  $\mathbf{A}\mathbf{x} \geq \mathbf{b}$  has the usual elementwise interpretation. Show that this linear programming problem can be converted to the problem in part b.

*Hint:* First consider diagonal  $\mathbf{F}_j$ .

- 22.21** Suppose that you have a cake and you need to divide it among  $n$  different children. Suppose that the  $i$ th child receives a fraction  $x_i$  of the cake. We will call the vector  $\mathbf{x} = [x_1, \dots, x_n]^\top$  an *allocation*. We require that every child receives at least some share of the cake, and that the entire cake is completely used up in the allocation. We also impose the additional condition that the first child ( $i = 1$ ) is allocated a share that is at least twice that of any other child. We say that the allocation is feasible if it meets all these requirements.

A feasible allocation  $\mathbf{x}$  is said to be *proportionally fair* if for any other allocation  $\mathbf{y}$ ,

$$\sum_{i=1}^n \frac{y_i - x_i}{x_i} \leq 0.$$

- a.** Let  $\Omega$  be the set of all feasible allocations. Show that  $\Omega$  is convex.
- b.** Show that a feasible allocation is proportionally fair if and only if it solves the following optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^n \log(x_i) \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

**22.22** Let  $U_i : \mathbb{R} \rightarrow \mathbb{R}$ ,  $U_i \in \mathcal{C}^1$ ,  $i = 1, \dots, n$ , be a set of concave increasing functions. Consider the optimization problem

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^n U_i(x_i) \\ & \text{subject to} && \sum_{i=1}^n x_i \leq C, \end{aligned}$$

where  $C > 0$  is a given constant.

- a. Show that the optimization problem above is a convex optimization problem.
- b. Show that  $\mathbf{x}^* = [x_1^*, \dots, x_n^*]^\top$  is an optimal solution to the optimization problem if and only if there exists a scalar  $\mu^* \geq 0$  such that  $x_i^* = \arg \max_x (U_i(x) - \mu^* x)$ . [The quantity  $U_i(x)$  has the interpretation of the “utility” of  $x$ , whereas  $\mu^*$  has the interpretation of a “price” per unit of  $x$ .]
- c. Show that  $\sum_{i=1}^n x_i^* = C$ .

**22.23** Give an example of a function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , a set  $\Omega = \{\mathbf{x} : g(\mathbf{x}) \leq 0\}$ , and a regular point  $\mathbf{x}^* \in \Omega$ , such that the following all hold simultaneously:

1.  $\mathbf{x}^*$  satisfies the FONC for set constraint  $\Omega$  (Theorem 6.1).
2.  $\mathbf{x}^*$  satisfies the KKT condition for inequality constraint  $g(\mathbf{x}) \leq 0$  (Theorem 21.1).
3.  $\mathbf{x}^*$  satisfies the SONC for set constraint  $\Omega$  (Theorem 6.2).
4.  $\mathbf{x}^*$  does not satisfy the SONC for inequality constraint  $g(\mathbf{x}) \leq 0$  (Theorem 21.2).

Be sure to show carefully that your choice of  $f$ ,  $\Omega = \{\mathbf{x} : g(\mathbf{x}) \leq 0\}$ , and  $\mathbf{x}^*$  satisfies all the conditions above simultaneously.

**22.24** This question is on duality theory for *nonlinear* programming problems, analogous to the theory for linear programming (Chapter 17). (A version for quadratic programming is considered in Exercise 17.24.)

Consider the following optimization problem:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex, each component of  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is convex, and  $f, \mathbf{g} \in \mathcal{C}^1$ . Let us call this problem the *primal* problem.

Define the *dual* of the problem above as

$$\begin{aligned} & \text{maximize } q(\boldsymbol{\mu}) \\ & \text{subject to } \boldsymbol{\mu} \geq \mathbf{0}, \end{aligned}$$

where  $q$  is defined by

$$q(\boldsymbol{\mu}) = \min_{\mathbf{x} \in \mathbb{R}^n} l(\mathbf{x}, \boldsymbol{\mu}),$$

with  $l(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x})$  the Lagrangian at  $\mathbf{x}, \boldsymbol{\mu}$ .

Prove the following results:

- a. If  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_0$  are feasible points in the primal and dual, respectively, then  $f(\mathbf{x}_0) \geq q(\boldsymbol{\mu}_0)$ . This is the *weak duality lemma* for nonlinear programming, analogous to Lemma 17.1.
- b. If  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_0$  are feasible points in the primal and dual, and  $f(\mathbf{x}_0) = q(\boldsymbol{\mu}_0)$ , then  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_0$  are optimal solutions to the primal and dual, respectively.
- c. If the primal has an optimal (feasible) solution, then so does the dual, and their objective function values are equal. (You may assume regularity.) This is the *duality theorem* for nonlinear programming, analogous to Theorem 17.2.

**22.25** Consider the matrix

$$\mathbf{M} = \begin{bmatrix} 1 & \gamma & -1 \\ \gamma & 1 & 2 \\ -1 & 2 & 5 \end{bmatrix},$$

where  $\gamma$  is a parameter.

- a. Find the Schur complement of  $\mathbf{M}(1,1)$ ;
- b. Find the Schur complement of  $\mathbf{M}(2:3, 2:3)$  (the bottom-right  $2 \times 2$  submatrix of  $\mathbf{M}$ , using MATLAB notation).

**22.26** Represent the Lyapunov inequality

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P} \mathbf{A} < 0,$$

where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & -2 \end{bmatrix},$$

as a canonical LMI.

**22.27** Let  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{R}$  be given matrices such that  $\mathbf{R} = \mathbf{R}^\top > 0$ . Suppose that we wish to find a symmetric positive definite matrix  $\mathbf{P}$  satisfying the following quadratic inequality:

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P}\mathbf{A} + \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^\top \mathbf{P} < 0.$$

Represent this inequality in the form of LMIs. (This inequality should not be confused with the *algebraic Riccati inequality*, which has a negative sign in front of the third term.)

**22.28** Let

$$\mathbf{A} = \begin{bmatrix} -0.9501 & -0.4860 & -0.4565 \\ -0.2311 & -0.8913 & -0.0185 \\ -0.6068 & -0.7621 & -0.8214 \end{bmatrix}.$$

Write a MATLAB program that finds a matrix  $\mathbf{P}$  satisfying  $0.1\mathbf{I}_3 \leq \mathbf{P} \leq \mathbf{I}_3$  and

$$\mathbf{A}^\top \mathbf{P} + \mathbf{P}\mathbf{A} \leq 0.$$

## CHAPTER 23

---

# ALGORITHMS FOR CONSTRAINED OPTIMIZATION

---

### 23.1 Introduction

In Part II we discussed algorithms for solving *unconstrained* optimization problems. In this chapter we present some simple algorithms for solving special *constrained* optimization problems. The methods here build on those of Part II.

We begin our presentation in the next section with a discussion of *projected methods*, including a treatment of projected gradient methods for problems with linear equality constraints. We then consider *Lagrangian methods*. Finally, we consider *penalty methods*. This chapter is intended as an introduction to ideas underlying methods for solving constrained optimization problems. For an in-depth coverage of the subject, we refer the reader to [11].

### 23.2 Projections

The optimization algorithms considered in Part II have the general form

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + \alpha_k \boldsymbol{d}^{(k)},$$

where  $\mathbf{d}^{(k)}$  is typically a function of  $\nabla f(\mathbf{x}^{(k)})$ . The value of  $\mathbf{x}^{(k)}$  is not constrained to lie inside any particular set. Such an algorithm is not immediately applicable to solving constrained optimization problems in which the decision variable is required to lie within a prespecified constraint set.

Consider the optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

If we use the algorithm above to solve this constrained problem, the iterates  $\mathbf{x}^{(k)}$  may not satisfy the constraints. Therefore, we need to modify the algorithms to take into account the presence of the constraints. A simple modification involves the introduction of a *projection*. The idea is as follows. If  $\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$  is in  $\Omega$ , then we set  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$  as usual. If, on the other hand,  $\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$  is not in  $\Omega$ , then we “project” it back into  $\Omega$  before setting  $\mathbf{x}^{(k+1)}$ .

To illustrate the projection method, consider the case where the constraint set  $\Omega \subset \mathbb{R}^n$  is given by

$$\Omega = \{\mathbf{x} : l_i \leq x_i \leq u_i, i = 1, \dots, n\}.$$

In this case,  $\Omega$  is a “box” in  $\mathbb{R}^n$ ; for this reason, this form of  $\Omega$  is called a *box constraint*. Given a point  $\mathbf{x} \in \mathbb{R}^n$ , define  $\mathbf{y} = \mathbf{\Pi}[\mathbf{x}] \in \mathbb{R}^n$  by

$$y_i = \min\{u_i, \max\{l_i, x_i\}\} = \begin{cases} u_i & \text{if } x_i > u_i \\ x_i & \text{if } l_i \leq x_i \leq u_i \\ l_i & \text{if } x_i < l_i. \end{cases}$$

The point  $\mathbf{\Pi}[\mathbf{x}]$  is called the *projection* of  $\mathbf{x}$  onto  $\Omega$ . Note that  $\mathbf{\Pi}[\mathbf{x}]$  is actually the “closest” point in  $\Omega$  to  $\mathbf{x}$ . Using the projection operator  $\mathbf{\Pi}$ , we can modify the previous unconstrained algorithm as follows:

$$\mathbf{x}^{(k+1)} = \mathbf{\Pi}[\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}].$$

Note that the iterates  $\mathbf{x}^{(k)}$  now all lie inside  $\Omega$ . We call the algorithm above a *projected algorithm*.

In the more general case, we can define the projection onto  $\Omega$ :

$$\mathbf{\Pi}[\mathbf{x}] = \arg \min_{\mathbf{z} \in \Omega} \|\mathbf{z} - \mathbf{x}\|.$$

In this case,  $\mathbf{\Pi}[\mathbf{x}]$  is again the “closest” point in  $\Omega$  to  $\mathbf{x}$ . This projection operator is well-defined only for certain types of constraint sets: for example, closed convex sets (see Exercise 22.19). For some sets  $\Omega$ , the “arg min” above is not well-defined. If the projection  $\mathbf{\Pi}$  is well-defined, we can similarly apply the projected algorithm

$$\mathbf{x}^{(k+1)} = \mathbf{\Pi}[\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}].$$

In some cases, there is a formula for computing  $\Pi[\mathbf{x}]$ . For example, if  $\Omega$  represents a box constraint as described above, then the formula given previously can be used. Another example is where  $\Omega$  is a linear variety, which is discussed in the next section. In general, even if the projection  $\Pi$  is well-defined, computation of  $\Pi[\mathbf{x}]$  for a given  $\mathbf{x}$  may not be easy. Often, the projection  $\Pi[\mathbf{x}]$  may have to be computed numerically. However, the numerical computation of  $\Pi[\mathbf{x}]$  itself entails solving an optimization algorithm. Indeed, the computation of  $\Pi[\mathbf{x}]$  may be as difficult as the original optimization problem, as is the case in the following example:

$$\begin{aligned} & \text{minimize} && \|\mathbf{x}\|^2 \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

Note that the solution to the problem in this case can be written as  $\Pi[\mathbf{0}]$ . Therefore, if  $\mathbf{0} \notin \Omega$ , the computation of a projection is equivalent to solving the given optimization problem.

As an example, consider the projection method applied specifically to the gradient algorithm (see Chapter 8). Recall that the vector  $-\nabla f(\mathbf{x})$  points in the direction of maximum rate of decrease of  $f$  at  $\mathbf{x}$ . This was the basis for gradient methods for unconstrained optimization, which have the form  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$ , where  $\alpha_k$  is the step size. The choice of the step size  $\alpha_k$  depends on the particular gradient algorithm. For example, recall that in the steepest descent algorithm,  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$ .

The projected version of the gradient algorithm has the form

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})].$$

We refer to the above as the *projected gradient algorithm*.

**Example 23.1** Consider the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \\ & \text{subject to} && \|\mathbf{x}\|^2 = 1, \end{aligned}$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . Suppose that we apply a *fixed-step-size projected gradient algorithm* to this problem.

- a. Derive a formula for the update equation for the algorithm (i.e., write down an explicit formula for  $\mathbf{x}^{(k+1)}$  as a function of  $\mathbf{x}^{(k)}$ ,  $\mathbf{Q}$ , and the fixed step size  $\alpha$ ). You may assume that the argument in the projection operator to obtain  $\mathbf{x}^{(k)}$  is never zero.
- b. Is it possible for the algorithm not to converge to an optimal solution even if the step size  $\alpha > 0$  is taken to be arbitrarily small?
- c. Show that for  $0 < \alpha < 1/\lambda_{\max}$  (where  $\lambda_{\max}$  is the largest eigenvalue of  $\mathbf{Q}$ ), the fixed-step-size projected gradient algorithm (with step size  $\alpha$ )

converges to an optimal solution, provided that  $\mathbf{x}^{(0)}$  is not orthogonal to the eigenvectors of  $\mathbf{Q}$  corresponding to the smallest eigenvalue. (Assume that the eigenvalues are distinct.)

*Solution:*

- a. The projection operator in this case simply maps any vector to the closest point on the unit circle. Therefore, the projection operator is given by  $\Pi[\mathbf{x}] = \mathbf{x}/\|\mathbf{x}\|$ , provided that  $\mathbf{x} \neq \mathbf{0}$ . The update equation is

$$\mathbf{x}^{(k+1)} = \beta_k(\mathbf{x}^{(k)} - \alpha\mathbf{Q}\mathbf{x}^{(k)}) = \beta_k(\mathbf{I} - \alpha\mathbf{Q})\mathbf{x}^{(k)},$$

where  $\beta_k = 1/\|(\mathbf{I} - \alpha\mathbf{Q})\mathbf{x}^{(k)}\|$  (i.e., it is whatever constant scaling is needed to make  $\mathbf{x}^{(k+1)}$  have unit norm).

- b. If we start with  $\mathbf{x}^{(0)}$  being an eigenvector of  $\mathbf{Q}$ , then  $\mathbf{x}^{(k)} = \mathbf{x}^{(0)}$  for all  $k$ . Therefore, if the corresponding eigenvalue is not the smallest, then clearly the algorithm is stuck at a point that is not optimal.

- c. We have

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \beta_k(\mathbf{I} - \alpha\mathbf{Q})\mathbf{x}^{(k)} \\ &= \beta_k(\mathbf{I} - \alpha\mathbf{Q})(y_1^{(k)}\mathbf{v}_1 + \cdots + y_n^{(k)}\mathbf{v}_n) \\ &= \beta_k(y_1^{(k)}(\mathbf{I} - \alpha\mathbf{Q})\mathbf{v}_1 + \cdots + y_n^{(k)}(\mathbf{I} - \alpha\mathbf{Q})\mathbf{v}_n).\end{aligned}$$

But  $(\mathbf{I} - \alpha\mathbf{Q})\mathbf{v}_i = (1 - \alpha\lambda_i)\mathbf{v}_i$ , where  $\lambda_i$  is the eigenvalue corresponding to  $\mathbf{v}_i$ . Hence,

$$\mathbf{x}^{(k+1)} = \beta_k(y_1^{(k)}(1 - \alpha\lambda_1)\mathbf{v}_1 + \cdots + y_n^{(k)}(1 - \alpha\lambda_n)\mathbf{v}_n),$$

which means that  $y_i^{(k+1)} = \beta_k y_i^{(k)}(1 - \alpha\lambda_i)$ . In other words,  $y_i^{(k)} = \beta^{(k)} y_i^{(0)}(1 - \alpha\lambda_i)^k$ , where  $\beta^{(k)} = \prod_{i=0}^{k-1} \beta_k$ . We rewrite  $\mathbf{x}^{(k)}$  as

$$\begin{aligned}\mathbf{x}^{(k)} &= \sum_{i=1}^n y_i^{(k)} \mathbf{v}_i \\ &= y_1^{(k)} \left( \mathbf{v}_1 + \sum_{i=2}^n \frac{y_i^{(k)}}{y_1^{(k)}} \mathbf{v}_i \right).\end{aligned}$$

Assuming that  $y_1^{(0)} \neq 0$ , we obtain

$$\frac{y_i^{(k)}}{y_1^{(k)}} = \frac{y_i^{(0)}(1 - \alpha\lambda_i)^k}{y_1^{(0)}(1 - \alpha\lambda_1)^k} = \frac{y_i^{(0)}}{y_1^{(0)}} \left( \frac{1 - \alpha\lambda_i}{1 - \alpha\lambda_1} \right)^k.$$

Using the fact that  $(1 - \alpha\lambda_i)/(1 - \alpha\lambda_1) < 1$  (because the  $\lambda_i > \lambda_1$  for  $i > 1$  and  $\alpha < 1/\lambda_{\max}$ ), we deduce that

$$\frac{y_i^{(k)}}{y_1^{(k)}} \rightarrow 0,$$

which implies that  $\mathbf{x}^{(k)} \rightarrow \mathbf{v}_1$ , as required. ■

### 23.3 Projected Gradient Methods with Linear Constraints

In this section we consider optimization problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ ,  $\mathbf{b} \in \mathbb{R}^m$ . We assume throughout that  $f \in \mathcal{C}^1$ . In the problem above, the constraint set is  $\Omega = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ . The specific structure of the constraint set allows us to compute the projection operator  $\Pi$  using the *orthogonal projector* (see Section 3.3). Specifically,  $\Pi[\mathbf{x}]$  can be defined using the orthogonal projector matrix  $\mathbf{P}$  given by

$$\mathbf{P} = \mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}$$

(see Example 12.5). Two important properties of the orthogonal projector  $\mathbf{P}$  that we use in this section are (see Theorem 3.5):

1.  $\mathbf{P} = \mathbf{P}^\top$ .
2.  $\mathbf{P}^2 = \mathbf{P}$ .

Another property of the orthogonal projector that we need in our discussion is given in the following lemma.

**Lemma 23.1** *Let  $\mathbf{v} \in \mathbb{R}^n$ . Then,  $\mathbf{P}\mathbf{v} = \mathbf{0}$  if and only if  $\mathbf{v} \in \mathcal{R}(\mathbf{A}^\top)$ . In other words,  $\mathcal{N}(\mathbf{P}) = \mathcal{R}(\mathbf{A}^\top)$ . Moreover,  $\mathbf{A}\mathbf{v} = \mathbf{0}$  if and only if  $\mathbf{v} \in \mathcal{R}(\mathbf{P})$ ; that is,  $\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{P})$ .  $\square$*

*Proof.*  $\Rightarrow$ : We have

$$\begin{aligned} \mathbf{P}\mathbf{v} &= (\mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A})\mathbf{v} \\ &= \mathbf{v} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{v}. \end{aligned}$$

If  $\mathbf{P}\mathbf{v} = \mathbf{0}$ , then

$$\mathbf{v} = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{v}$$

and hence  $\mathbf{v} \in \mathcal{R}(\mathbf{A}^\top)$ .

$\Leftarrow$ : Suppose that there exists  $\mathbf{u} \in \mathbb{R}^m$  such that  $\mathbf{v} = \mathbf{A}^\top \mathbf{u}$ . Then,

$$\begin{aligned} \mathbf{P}\mathbf{v} &= (\mathbf{I}_n - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A})\mathbf{A}^\top \mathbf{u} \\ &= \mathbf{A}^\top \mathbf{u} - \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{A}^\top \mathbf{u} \\ &= \mathbf{0}. \end{aligned}$$

Hence, we have proved that  $\mathcal{N}(\mathbf{P}) = \mathcal{R}(\mathbf{A}^\top)$ .

Using an argument similar to that above, we can show that  $\mathcal{N}(\mathbf{A}) = \mathcal{R}(\mathbf{P})$ .  $\blacksquare$

Recall that in unconstrained optimization, the first-order necessary condition for a point  $\mathbf{x}^*$  to be a local minimizer is  $\nabla f(\mathbf{x}^*) = \mathbf{0}$  (see Section 6.2). In optimization problems with equality constraints, the Lagrange condition plays the role of the first-order necessary condition (see Section 20.4). When the constraint set takes the form  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ , the Lagrange condition can be written as  $\mathbf{P}\nabla f(\mathbf{x}^*) = \mathbf{0}$ , as stated in the following proposition.

**Proposition 23.1** *Let  $\mathbf{x}^* \in \mathbb{R}^n$  be a feasible point. Then,  $\mathbf{P}\nabla f(\mathbf{x}^*) = \mathbf{0}$  if and only if  $\mathbf{x}^*$  satisfies the Lagrange condition.  $\square$*

*Proof.* By Lemma 23.1,  $\mathbf{P}\nabla f(\mathbf{x}^*) = \mathbf{0}$  if and only if we have  $\nabla f(\mathbf{x}^*) \in \mathcal{R}(\mathbf{A}^\top)$ . This is equivalent to the condition that there exists  $\boldsymbol{\lambda}^* \in \mathbb{R}^m$  such that  $\nabla f(\mathbf{x}^*) + \mathbf{A}^\top \boldsymbol{\lambda}^* = \mathbf{0}$ , which together with the feasibility equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , constitutes the Lagrange condition.  $\blacksquare$

Recall that the projected gradient algorithm has the form

$$\mathbf{x}^{(k+1)} = \mathbf{\Pi}[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})].$$

For the case where the constraints are linear, it turns out that we can express the projection  $\mathbf{\Pi}$  in terms of the matrix  $\mathbf{P}$  as follows:

$$\mathbf{\Pi}[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})] = \mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)}),$$

assuming that  $\mathbf{x}^{(k)} \in \Omega$ . Although the formula above can be derived algebraically (see Exercise 23.4), it is more insightful to derive the formula using a geometric argument, as follows. In our constrained optimization problem, the vector  $-\nabla f(\mathbf{x})$  is not necessarily a feasible direction. In other words, if  $\mathbf{x}^{(k)}$  is a feasible point and we apply the algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$ , then  $\mathbf{x}^{(k+1)}$  need not be feasible. This problem can be overcome by replacing  $-\nabla f(\mathbf{x}^{(k)})$  by a vector that points in a feasible direction. Note that the set of feasible directions is simply the nullspace  $\mathcal{N}(\mathbf{A})$  of the matrix  $\mathbf{A}$ . Therefore, we should first project the vector  $-\nabla f(\mathbf{x})$  onto  $\mathcal{N}(\mathbf{A})$ . This projection is equivalent to multiplication by the matrix  $\mathbf{P}$ . In summary, in the projection gradient algorithm, we update  $\mathbf{x}^{(k)}$  according to the equation

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)}).$$

The projected gradient algorithm has the following property.

**Proposition 23.2** *In a projected gradient algorithm, if  $\mathbf{x}^{(0)}$  is feasible, then each  $\mathbf{x}^{(k)}$  is feasible; that is, for each  $k \geq 0$ ,  $\mathbf{A}\mathbf{x}^{(k)} = \mathbf{b}$ .  $\square$*

*Proof.* We proceed by induction. The result holds for  $k = 0$  by assumption. Suppose now that  $\mathbf{A}\mathbf{x}^{(k)} = \mathbf{b}$ . We now show that  $\mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{b}$ . To show this, first observe that  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \in \mathcal{N}(\mathbf{A})$ . Therefore,

$$\begin{aligned} \mathbf{A}\mathbf{x}^{(k+1)} &= \mathbf{A}(\mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)})) \\ &= \mathbf{A}\mathbf{x}^{(k)} - \alpha_k \mathbf{A}\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \\ &= \mathbf{b}, \end{aligned}$$

which completes the proof. ■

The projected gradient algorithm updates  $\mathbf{x}^{(k)}$  in the direction of  $-\mathbf{P}\nabla f(\mathbf{x}^{(k)})$ . This vector points in the direction of maximum rate of decrease of  $f$  at  $\mathbf{x}^{(k)}$  along the surface defined by  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , as described in the following argument. Let  $\mathbf{x}$  be any feasible point and  $\mathbf{d}$  a feasible direction such that  $\|\mathbf{d}\| = 1$ . The rate of increase of  $f$  at  $\mathbf{x}$  in the direction  $\mathbf{d}$  is  $\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle$ . Next, we note that because  $\mathbf{d}$  is a feasible direction, it lies in  $\mathcal{N}(\mathbf{A})$  and hence by Lemma 23.1, we have  $\mathbf{d} \in \mathcal{R}(\mathbf{P}) = \mathcal{R}(\mathbf{P}^\top)$ . So, there exists  $\mathbf{v}$  such that  $\mathbf{d} = \mathbf{P}\mathbf{v}$ . Hence,

$$\langle \nabla f(\mathbf{x}), \mathbf{d} \rangle = \langle \nabla f(\mathbf{x}), \mathbf{P}^\top \mathbf{v} \rangle = \langle \mathbf{P}\nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

By the Cauchy-Schwarz inequality,

$$\langle \mathbf{P}\nabla f(\mathbf{x}), \mathbf{v} \rangle \leq \|\mathbf{P}\nabla f(\mathbf{x})\| \|\mathbf{v}\|$$

with equality if and only if the direction of  $\mathbf{v}$  is parallel with the direction of  $\mathbf{P}\nabla f(\mathbf{x})$ . Therefore, the vector  $-\mathbf{P}\nabla f(\mathbf{x})$  points in the direction of maximum rate of decrease of  $f$  at  $\mathbf{x}$  among all feasible directions.

Following the discussion in Chapter 8 for gradient methods in unconstrained optimization, we suggest the following gradient method for our constrained problem. Suppose that we have a starting point  $\mathbf{x}^{(0)}$ , which we assume is feasible; that is,  $\mathbf{A}\mathbf{x}^{(0)} = \mathbf{b}$ . Consider the point  $\mathbf{x} = \mathbf{x}^{(0)} - \alpha \mathbf{P}\nabla f(\mathbf{x}^{(0)})$ , where  $\alpha \in \mathbb{R}$ . As usual, the scalar  $\alpha$  is called the step size. By the discussion above,  $\mathbf{x}$  is also a feasible point. Using a Taylor series expansion of  $f$  about  $\mathbf{x}^{(0)}$  and the fact that  $\mathbf{P} = \mathbf{P}^2 = \mathbf{P}^\top \mathbf{P}$ , we get

$$\begin{aligned} f(\mathbf{x}^{(0)} - \alpha \mathbf{P}\nabla f(\mathbf{x}^{(0)})) &= f(\mathbf{x}^{(0)}) - \alpha \nabla f(\mathbf{x}^{(0)})^\top \mathbf{P}\nabla f(\mathbf{x}^{(0)}) + o(\alpha) \\ &= f(\mathbf{x}^{(0)}) - \alpha \|\mathbf{P}\nabla f(\mathbf{x}^{(0)})\|^2 + o(\alpha). \end{aligned}$$

Thus, if  $\mathbf{P}\nabla f(\mathbf{x}^{(0)}) \neq 0$ , that is,  $\mathbf{x}^{(0)}$  does not satisfy the Lagrange condition, then we can choose an  $\alpha$  sufficiently small such that  $f(\mathbf{x}) < f(\mathbf{x}^{(0)})$ , which means that  $\mathbf{x} = \mathbf{x}^{(0)} - \alpha \mathbf{P}\nabla f(\mathbf{x}^{(0)})$  is an improvement over  $\mathbf{x}^{(0)}$ . This is the basis for the projected gradient algorithm  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)})$ , where the initial point  $\mathbf{x}^{(0)}$  satisfies  $\mathbf{A}\mathbf{x}^{(0)} = \mathbf{b}$  and  $\alpha_k$  is some step size. As for unconstrained gradient methods, the choice of  $\alpha_k$  determines the behavior of the algorithm. For small step sizes, the algorithm progresses slowly, while large step sizes may result in a zigzagging path. A well-known variant of the projected gradient algorithm is the *projected steepest descent algorithm*, where  $\alpha_k$  is given by

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{P}\nabla f(\mathbf{x}^{(k)})).$$

The following theorem states that the projected steepest descent algorithm is a descent algorithm, in the sense that at each step the value of the objective function decreases.

**Theorem 23.1** If  $\{\mathbf{x}^{(k)}\}$  is the sequence of points generated by the projected steepest descent algorithm and if  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , then  $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ .  $\square$

*Proof.* First, recall that

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{P}\nabla f(\mathbf{x}^{(k)}),$$

where  $\alpha_k \geq 0$  is the minimizer of

$$\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{P}\nabla f(\mathbf{x}^{(k)}))$$

over all  $\alpha \geq 0$ . Thus, for  $\alpha \geq 0$ , we have

$$\phi_k(\alpha_k) \leq \phi_k(\alpha).$$

By the chain rule,

$$\begin{aligned}\phi'_k(0) &= \frac{d\phi_k}{d\alpha}(0) \\ &= -\nabla f(\mathbf{x}^{(k)} - 0 \mathbf{P}\nabla f(\mathbf{x}^{(k)}))^T \mathbf{P}\nabla f(\mathbf{x}^{(k)}) \\ &= -\nabla f(\mathbf{x}^{(k)})^T \mathbf{P}\nabla f(\mathbf{x}^{(k)}).\end{aligned}$$

Using the fact that  $\mathbf{P} = \mathbf{P}^2 = \mathbf{P}^\top \mathbf{P}$ , we get

$$\phi'_k(0) = -\nabla f(\mathbf{x}^{(k)})^T \mathbf{P}^\top \mathbf{P}\nabla f(\mathbf{x}^{(k)}) = -\|\mathbf{P}\nabla f(\mathbf{x}^{(k)})\|^2 < 0,$$

because  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$  by assumption. Thus, there exists  $\bar{\alpha} > 0$  such that  $\phi_k(0) > \phi_k(\alpha)$  for all  $\alpha \in (0, \bar{\alpha}]$ . Hence,

$$f(\mathbf{x}^{(k+1)}) = \phi_k(\alpha_k) \leq \phi_k(\bar{\alpha}) < \phi_k(0) = f(\mathbf{x}^{(k)}),$$

which completes the proof of the theorem.  $\blacksquare$

In Theorem 23.1 we needed the assumption that  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$  to prove that the algorithm possesses the descent property. If for some  $k$ , we have  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$ , then by Proposition 23.1 the point  $\mathbf{x}^{(k)}$  satisfies the Lagrange condition. This condition can be used as a stopping criterion for the algorithm. Note that in this case,  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ . For the case where  $f$  is a convex function, the condition  $\mathbf{P}\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$  is, in fact, equivalent to  $\mathbf{x}^{(k)}$  being a global minimizer of  $f$  over the constraint set  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ . We show this in the following proposition.

**Proposition 23.3** The point  $\mathbf{x}^* \in \mathbb{R}^n$  is a global minimizer of a convex function  $f$  over  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$  if and only if  $\mathbf{P}\nabla f(\mathbf{x}^*) = \mathbf{0}$ .  $\square$

*Proof.* We first write  $\mathbf{h}(\mathbf{x}) = \mathbf{A}\mathbf{x} - \mathbf{b}$ . Then, the constraints can be written as  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ , and the problem is of the form considered in earlier chapters.

Note that  $D\mathbf{h}(\mathbf{x}) = \mathbf{A}$ . Hence,  $\mathbf{x}^* \in \mathbb{R}^n$  is a global minimizer of  $f$  if and only if the Lagrange condition holds (see Theorem 22.8). By Proposition 23.1, this is true if and only if  $P\nabla f(\mathbf{x}^*) = \mathbf{0}$ , and this completes the proof. ■

For an application of the projected steepest descent algorithm to minimum fuel and minimum amplitude control problems in linear discrete systems, see [78].

## 23.4 Lagrangian Algorithms

In this section we consider an optimization method based on the Lagrangian function (see Section 20.4). The basic idea is to use gradient algorithms to update simultaneously the decision variable and Lagrange multiplier vector. We consider first the case with equality constraints, followed by inequality constraints.

### Lagrangian Algorithm for Equality Constraints

Consider the following optimization problem with equality constraints:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0} \end{aligned}$$

where  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . Recall that for this problem the Lagrangian function is given by

$$l(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{h}(\mathbf{x}).$$

Assume that  $f, \mathbf{h} \in \mathcal{C}^2$ ; as usual, denote the Hessian of the Lagrangian by  $\mathbf{L}(\mathbf{x}, \boldsymbol{\lambda})$ .

The Lagrangian algorithm for this problem is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + D\mathbf{h}(\mathbf{x}^{(k)})^\top \boldsymbol{\lambda}^{(k)}), \\ \boldsymbol{\lambda}^{(k+1)} &= \boldsymbol{\lambda}^{(k)} + \beta_k \mathbf{h}(\mathbf{x}^{(k)}). \end{aligned}$$

Notice that the update equation for  $\mathbf{x}^{(k)}$  is a gradient algorithm for minimizing the Lagrangian with respect to its  $\mathbf{x}$  argument, and the update equation for  $\boldsymbol{\lambda}^{(k)}$  is a gradient algorithm for maximizing the Lagrangian with respect to its  $\boldsymbol{\lambda}$  argument. Because only the gradient is used, the method is also called the *first-order Lagrangian algorithm*.

The following lemma establishes that if the algorithm converges, the limit must satisfy the Lagrange condition. More specifically, the lemma states that any *fixed point* of the algorithm must satisfy the Lagrange condition. A fixed point of an update algorithm is simply a point with the property that when updated using the algorithm, the resulting point is equal to the given point. For the case of the Lagrangian algorithm, which updates both  $\mathbf{x}^{(k)}$  and

$\lambda^{(k)}$  vectors, a fixed point is a *pair* of vectors. If the Lagrangian algorithm converges, the limit must be a fixed point. We omit the proof of the lemma because it follows easily by inspection.

**Lemma 23.2** *For the Lagrangian algorithm for updating  $\mathbf{x}^{(k)}$  and  $\boldsymbol{\lambda}^{(k)}$ , the pair  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  is a fixed point if and only if it satisfies the Lagrange condition.*

□

Below, we use  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  to denote a pair satisfying the Lagrange condition. Assume that  $L(\mathbf{x}^*, \boldsymbol{\lambda}^*) > 0$ . Also assume that  $\mathbf{x}^*$  is a *regular* point. With these assumptions, we are now ready to state and prove that the algorithm is locally convergent. For simplicity, we will take  $\alpha_k$  and  $\beta_k$  to be fixed constants (not depending on  $k$ ), denoted  $\alpha$  and  $\beta$ , respectively.

**Theorem 23.2** *For the Lagrangian algorithm for updating  $\mathbf{x}^{(k)}$  and  $\boldsymbol{\lambda}^{(k)}$ , provided that  $\alpha$  and  $\beta$  are sufficiently small, there is a neighborhood of  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  such that if the pair  $(\mathbf{x}^{(0)}, \boldsymbol{\lambda}^{(0)})$  is in this neighborhood, then the the algorithm converges to  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  with at least a linear order of convergence.*

□

*Proof.* We can rescale  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  by appropriate constants (so that the assumptions are preserved) and effectively change the relative values of the step sizes for the update equations. Therefore, without loss of generality, we can take  $\beta = \alpha$ .

We will set up our proof by introducing some convenient notation. Given a pair  $(\mathbf{x}, \boldsymbol{\lambda})$ , let  $\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top$  be the  $(n+m)$ -vector constructed by concatenating  $\mathbf{x}$  and  $\boldsymbol{\lambda}$ . Similarly define  $\mathbf{w}^{(k)} = [\mathbf{x}^{(k)\top}, \boldsymbol{\lambda}^{(k)\top}]^\top$  and  $\mathbf{w}^* = [\mathbf{x}^{*\top}, \boldsymbol{\lambda}^{*\top}]^\top$ . Define the map  $\mathbf{U} : \mathbb{R}^{n+m} \rightarrow \mathbb{R}^{n+m}$  by

$$\mathbf{U}(\mathbf{w}) = \begin{bmatrix} \mathbf{x} - \alpha(\nabla f(\mathbf{x}) + D\mathbf{h}(\mathbf{x})^\top \boldsymbol{\lambda}) \\ \boldsymbol{\lambda} + \alpha\mathbf{h}(\mathbf{x}) \end{bmatrix}.$$

Then, the Lagrangian algorithm can be rewritten as

$$\mathbf{w}^{(k+1)} = \mathbf{U}(\mathbf{w}^{(k)}).$$

We now write  $\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\|$  in terms of  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\|$ , where  $\|\cdot\|$  denotes the usual Euclidean norm. By Lemma 23.2,  $\mathbf{w}^* = [\mathbf{x}^{*\top}, \boldsymbol{\lambda}^{*\top}]^\top$  is a fixed point of  $\mathbf{w}^{(k+1)} = \mathbf{U}(\mathbf{w}^{(k)})$ . Therefore,

$$\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| = \|\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*)\|.$$

Let  $D\mathbf{U}$  be the (matrix) derivative of  $\mathbf{U}$ :

$$D\mathbf{U}(\mathbf{w}) = \mathbf{I} + \alpha \begin{bmatrix} -L(\mathbf{x}, \boldsymbol{\lambda}) & -D\mathbf{h}(\mathbf{x})^\top \\ D\mathbf{h}(\mathbf{x}) & \mathbf{O} \end{bmatrix}.$$

By the mean value theorem (see Theorem 5.9),

$$\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*) = \mathbf{G}(\mathbf{w}^{(k)})(\mathbf{w}^{(k)} - \mathbf{w}^*),$$

where  $\mathbf{G}(\mathbf{w}^{(k)})$  is a matrix whose rows are the rows of  $D\mathbf{U}$  evaluated at points that lie on the line segment joining  $\mathbf{w}^{(k)}$  and  $\mathbf{w}^*$  (these points may differ from row to row). Taking norms of both sides of the equation above,

$$\|\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*)\| \leq \|\mathbf{G}(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\|.$$

Finally, combining the above, we have

$$\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \|\mathbf{G}(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\|.$$

We now claim that for sufficiently small  $\alpha > 0$ ,  $\|D\mathbf{U}(\mathbf{w}^*)\| < 1$ . Our argument here follows [11, Section 4.4]. Let

$$\mathbf{M} = \begin{bmatrix} -\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) & -D\mathbf{h}(\mathbf{x}^*)^\top \\ D\mathbf{h}(\mathbf{x}^*) & \mathbf{O} \end{bmatrix},$$

so that  $D\mathbf{U}(\mathbf{w}^*) = \mathbf{I} + \alpha\mathbf{M}$ . Hence, to prove the claim, it suffices to show that the eigenvalues of  $\mathbf{M}$  all lie in the open left-half complex plane.

For any complex vector  $\mathbf{y}$ , let  $\mathbf{y}^H$  represent its complex conjugate transpose (or Hermitian) and  $\Re(\mathbf{y})$  its real part. Let  $\lambda$  be an eigenvalue of  $\mathbf{M}$  and  $\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top \neq \mathbf{0}$  be a corresponding eigenvector. Now,  $\Re(\mathbf{w}^H \mathbf{M} \mathbf{w}) = \Re(\lambda) \|\mathbf{w}\|^2$ . However, from the structure of  $\mathbf{M}$ , we can readily see that

$$\begin{aligned} \Re(\mathbf{w}^H \mathbf{M} \mathbf{w}) &= -\Re(\mathbf{x}^H \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}) - \Re(\mathbf{x}^H D\mathbf{h}(\mathbf{x}^*)^\top \boldsymbol{\lambda}) + \Re(\boldsymbol{\lambda}^H D\mathbf{h}(\mathbf{x}^*) \mathbf{x}) \\ &= -\Re(\mathbf{x}^H \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}). \end{aligned}$$

By the assumption that  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) > 0$ , we know that  $\Re(\mathbf{x}^H \mathbf{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}) > 0$  if  $\mathbf{x} \neq \mathbf{0}$ . Therefore, comparing the two equations above, we deduce that  $\Re(\lambda) < 0$ , as required, provided that  $\mathbf{x}$  is nonzero, as we now demonstrate.

Now, suppose that  $\mathbf{x} = \mathbf{0}$ . Because  $\mathbf{w}$  is an eigenvector of  $\mathbf{M}$ , we have  $\mathbf{M}\mathbf{w} = \lambda\mathbf{w}$ . Extracting the first  $n$  components, we have  $D\mathbf{h}(\mathbf{x}^*)^\top \boldsymbol{\lambda} = \mathbf{0}$ . By the regularity assumption, we deduce that  $\boldsymbol{\lambda} = \mathbf{0}$ . This contradicts the assumption that  $\mathbf{w} \neq \mathbf{0}$ . Hence we conclude that  $\mathbf{x} \neq \mathbf{0}$ , which completes the proof of our claim that for sufficiently small  $\alpha > 0$ ,  $\|D\mathbf{U}(\mathbf{w}^*)\| < 1$ .

The result of the foregoing claim allows us to pick constants  $\eta > 0$  and  $\kappa < 1$  such that for all  $\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\lambda}^\top]^\top$  satisfying  $\|\mathbf{w} - \mathbf{w}^*\| \leq \eta$ , we have  $\|\mathbf{G}(\mathbf{w})\| \leq \kappa$  (this follows from the continuity of  $D\mathbf{U}$  and norms).

To complete the proof, suppose that  $\|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \eta$ . We will show by induction that for all  $k \geq 0$ ,  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \eta$  and  $\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \kappa \|\mathbf{w}^{(k)} - \mathbf{w}^*\|$ , from which we conclude that  $\mathbf{w}^{(k)}$  converges to  $\mathbf{w}^*$  with at least linear order of convergence. For  $k = 0$ , the result follows because  $\|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \eta$  by assumption, and

$$\|\mathbf{w}^{(1)} - \mathbf{w}^*\| \leq \|\mathbf{G}(\mathbf{w}^{(0)})\| \|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \kappa \|\mathbf{w}^{(0)} - \mathbf{w}^*\|,$$

which follows from  $\|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \eta$ . So suppose that the result holds for  $k$ . This implies that  $\|\mathbf{G}(\mathbf{w}^{(k)})\| \leq \kappa$ . To show the  $k + 1$  case, we write

$$\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \|\mathbf{G}(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \kappa \|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \eta.$$

This means that  $\|\mathbf{G}(\mathbf{w}^{(k+1)})\| \leq \kappa$ , from which we can write

$$\|\mathbf{w}^{(k+2)} - \mathbf{w}^*\| \leq \|\mathbf{G}(\mathbf{w}^{(k+1)})\| \|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \kappa \|\mathbf{w}^{(k+1)} - \mathbf{w}^*\|.$$

This completes the proof. ■

### Lagrangian Algorithm for Inequality Constraints

Consider the following optimization problem with inequality constraints:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \end{aligned}$$

where  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ . Recall that for this problem the Lagrangian function is given by

$$l(\mathbf{x}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}).$$

As before, assume that  $f, \mathbf{g} \in \mathcal{C}^2$ ; as usual, denote the Hessian of the Lagrangian by  $\mathbf{L}(\mathbf{x}, \boldsymbol{\mu})$ .

The Lagrangian algorithm for this problem is given by

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \alpha_k (\nabla f(\mathbf{x}^{(k)}) + D\mathbf{g}(\mathbf{x}^{(k)})^\top \boldsymbol{\mu}^{(k)}), \\ \boldsymbol{\mu}^{(k+1)} &= [\boldsymbol{\mu}^{(k)} + \beta_k \mathbf{g}(\mathbf{x}^{(k)})]_+, \end{aligned}$$

where  $[\cdot]_+ = \max\{\cdot, 0\}$  (applied componentwise). Notice that, as before, the update equation for  $\mathbf{x}^{(k)}$  is a gradient algorithm for minimizing the Lagrangian with respect to its  $\mathbf{x}$  argument. The update equation for  $\boldsymbol{\mu}^{(k)}$  is a *projected gradient algorithm* for maximizing the Lagrangian with respect to its  $\boldsymbol{\mu}$  argument. The reason for the projection is that the KKT multiplier vector is required to be nonnegative to satisfy the KKT condition.

The following lemma establishes that if the algorithm converges, the limit must satisfy the KKT condition. As before, we use the notion of a fixed point to state the result formally. The proof is omitted because the result follows easily by inspection.

**Lemma 23.3** *For the Lagrangian algorithm for updating  $\mathbf{x}^{(k)}$  and  $\boldsymbol{\mu}^{(k)}$ , the pair  $(\mathbf{x}^*, \boldsymbol{\mu}^*)$  is a fixed point if and only if it satisfies the KKT condition.*  $\square$

As before, we use the notation  $(\mathbf{x}^*, \boldsymbol{\mu}^*)$  to denote a pair satisfying the KKT condition. Assume that  $\mathbf{L}(\mathbf{x}^*, \boldsymbol{\mu}^*) > 0$ . Also assume that  $\mathbf{x}^*$  is a *regular* point. With these assumptions, we are now ready to state and prove that the algorithm is locally convergent. As before, we will take  $\alpha_k$  and  $\beta_k$  to

be fixed constants (not depending on  $k$ ), denoted  $\alpha$  and  $\beta$ , respectively. Our analysis examines the behavior of the algorithm in two phases. In the first phase, the “nonactive” multipliers decrease to zero in finite time and remain at zero thereafter. In the second phase, the  $\mathbf{x}^{(k)}$  iterates and the “active” multipliers converge jointly to their respective solutions, with at least a linear order of convergence.

**Theorem 23.3** *For the Lagrangian algorithm for updating  $\mathbf{x}^{(k)}$  and  $\boldsymbol{\mu}^{(k)}$ , provided that  $\alpha$  and  $\beta$  are sufficiently small, there is a neighborhood of  $(\mathbf{x}^*, \boldsymbol{\mu}^*)$  such that if the pair  $(\mathbf{x}^{(0)}, \boldsymbol{\mu}^{(0)})$  is in this neighborhood, then (1) the nonactive multipliers reduce to zero in finite time and remain at zero thereafter and (2) the algorithm converges to  $(\mathbf{x}^*, \boldsymbol{\mu}^*)$  with at least a linear order of convergence.*

□

*Proof.* As in the proof of Theorem 23.2, we can rescale  $\mathbf{x}$  and  $\boldsymbol{\mu}$  by appropriate constants (so that the assumptions are preserved) and effectively change the relative values of the step sizes for the update equations. Therefore, without loss of generality, we can take  $\beta = \alpha$ .

We set up our proof using the same vector notation as before. Given a pair  $(\mathbf{x}, \boldsymbol{\mu})$ , let  $\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\mu}^\top]^\top$  be the  $(n + p)$ -vector constructed by concatenating  $\mathbf{x}$  and  $\boldsymbol{\mu}$ . Similarly define  $\mathbf{w}^{(k)} = [\mathbf{x}^{(k)\top}, \boldsymbol{\mu}^{(k)\top}]^\top$  and  $\mathbf{w}^* = [\mathbf{x}^{*\top}, \boldsymbol{\mu}^{*\top}]^\top$ . Define the map  $\mathbf{U}$  as

$$\mathbf{U}(\mathbf{w}) = \begin{bmatrix} \mathbf{x} - \alpha(\nabla f(\mathbf{x}) + D\mathbf{g}(\mathbf{x})^\top \boldsymbol{\mu}) \\ \boldsymbol{\mu} + \alpha\mathbf{g}(\mathbf{x}) \end{bmatrix}.$$

Also, define the map  $\mathbf{\Pi}$  by

$$\mathbf{\Pi}[\mathbf{w}] = \begin{bmatrix} \mathbf{x} \\ [\boldsymbol{\mu}]_+ \end{bmatrix}.$$

Then, the update equations can be rewritten as

$$\mathbf{w}^{(k+1)} = \mathbf{\Pi}[\mathbf{U}(\mathbf{w}^{(k)})].$$

Because  $\mathbf{\Pi}$  is a projection onto the convex set  $\{\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\mu}^\top]^\top : \boldsymbol{\mu} \geq \mathbf{0}\}$ , it is a nonexpansive map (see [12, Proposition 3.2]), which means that  $\|\mathbf{\Pi}[\mathbf{v}] - \mathbf{\Pi}[\mathbf{w}]\| \leq \|\mathbf{v} - \mathbf{w}\|$ .

We now write  $\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\|$  in terms of  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\|$ , where  $\|\cdot\|$  denotes the usual Euclidean norm. By Lemma 23.3,  $\mathbf{w}^* = [\mathbf{x}^{*\top}, \boldsymbol{\mu}^{*\top}]^\top$  is a fixed point of  $\mathbf{w}^{(k+1)} = \mathbf{U}(\mathbf{w}^{(k)})$ . Therefore,

$$\begin{aligned} \|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| &= \|\mathbf{\Pi}[\mathbf{U}(\mathbf{w}^{(k)})] - \mathbf{\Pi}[\mathbf{U}(\mathbf{w}^*)]\| \\ &\leq \|\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*)\| \end{aligned}$$

by the nonexpansiveness of  $\Pi$ . Let  $D\mathbf{U}$  be the (matrix) derivative of  $\mathbf{U}$ :

$$D\mathbf{U}(\mathbf{w}) = \mathbf{I} + \alpha \begin{bmatrix} -\mathbf{L}(\mathbf{x}, \boldsymbol{\mu}) & -D\mathbf{g}(\mathbf{x})^\top \\ D\mathbf{g}(\mathbf{x}) & \mathbf{O} \end{bmatrix}.$$

By the mean value theorem,

$$\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*) = \mathbf{G}(\mathbf{w}^{(k)})(\mathbf{w}^{(k)} - \mathbf{w}^*),$$

where  $\mathbf{G}(\mathbf{w}^{(k)})$  is a matrix whose rows are the rows of  $D\mathbf{U}$  evaluated at points that lie on the line segment joining  $\mathbf{w}^{(k)}$  and  $\mathbf{w}^*$  (these points may differ from row to row). Taking norms of both sides of the equation above yields

$$\|\mathbf{U}(\mathbf{w}^{(k)}) - \mathbf{U}(\mathbf{w}^*)\| \leq \|\mathbf{G}(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\|.$$

Finally, combining the above, we obtain

$$\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \|\mathbf{G}(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\|.$$

Let  $\mathbf{g}_A$  represent those rows of  $\mathbf{g}$  corresponding to active constraints (at  $\mathbf{x}^*$ ) and  $\mathbf{g}_{\bar{A}}$  represent the remaining rows of  $\mathbf{g}$ . [Recall that by regularity,  $D\mathbf{g}_A(\mathbf{x}^*)$  has full rank.] Given a vector  $\boldsymbol{\mu}$ , we divide it into two subvectors  $\boldsymbol{\mu}_A$  and  $\boldsymbol{\mu}_{\bar{A}}$ , according to active and nonactive components, respectively. (Note that  $\boldsymbol{\mu}_{\bar{A}}^* = \mathbf{0}$ , the zero vector.) Next, write  $\mathbf{w}_A = [\mathbf{x}^\top, \boldsymbol{\mu}_A^\top]^\top$  and

$$\mathbf{U}_A(\mathbf{w}_A) = \begin{bmatrix} \mathbf{x} - \alpha(\nabla f(\mathbf{x}) + D\mathbf{g}_A(\mathbf{x})^\top \boldsymbol{\mu}_A) \\ \boldsymbol{\mu}_A + \alpha\mathbf{g}_A(\mathbf{x}) \end{bmatrix},$$

so that

$$D\mathbf{U}_A(\mathbf{w}_A) = \mathbf{I} + \alpha \begin{bmatrix} -\mathbf{L}(\mathbf{x}, \boldsymbol{\mu}_A) & -D\mathbf{g}_A(\mathbf{x})^\top \\ D\mathbf{g}_A(\mathbf{x}) & \mathbf{O} \end{bmatrix}.$$

Finally, let  $\mathbf{G}_A$  be such that  $\mathbf{U}_A(\mathbf{w}_A^{(k)}) - \mathbf{U}_A(\mathbf{w}_A^*) = \mathbf{G}_A(\mathbf{w}_A^{(k)})(\mathbf{w}_A^{(k)} - \mathbf{w}_A^*)$  (by the mean value theorem as before).

We organize the remainder of our proof into four claims.

**Claim 1:** For sufficiently small  $\alpha > 0$ ,  $\|D\mathbf{U}_A(\mathbf{w}_A^*)\| < 1$ .

The argument here parallels that of the proof of Theorem 23.2. So for the sake of brevity we omit the details.

The result of claim 1 allows us to pick constants  $\eta > 0$ ,  $\delta > 0$ , and  $\kappa_A < 1$  such that for all  $\mathbf{w} = [\mathbf{x}^\top, \boldsymbol{\mu}^\top]^\top$  satisfying  $\|\mathbf{w} - \mathbf{w}^*\| \leq \eta$ ,  $\|\mathbf{G}_A(\mathbf{w}_A)\| \leq \kappa_A$ , and  $\mathbf{g}_{\bar{A}}(\mathbf{x}) \leq -\delta\mathbf{e}$ , where  $\mathbf{e}$  is the vector with all components equal to 1. The first inequality follows from claim 1 and the continuity of  $D\mathbf{U}_A(\cdot)$  and  $\|\cdot\|$ . The second follows from the fact that the components of  $\mathbf{g}_{\bar{A}}(\mathbf{x}^*)$  are negative.

Let  $\kappa = \max\{\|\mathbf{G}(\mathbf{w})\| : \|\mathbf{w} - \mathbf{w}^*\| \leq \eta\}$ , which we assume to be at least 1; otherwise, set  $\kappa = 1$ . Now pick  $\varepsilon > 0$  such that  $\varepsilon\kappa^{\varepsilon/(\alpha\delta)} \leq \eta$ . We can do this because the left side of this inequality goes to 0 as  $\varepsilon \rightarrow 0$ . Assume for

convenience that  $k_0 = \varepsilon/(\alpha\delta)$  is an integer; otherwise, replace all instances of  $\varepsilon/(\alpha\delta)$  by the smallest integer that exceeds it (i.e., round it up to the closest integer).

For the remainder of this proof, let  $\mathbf{w}^{(0)}$  satisfy  $\|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \varepsilon$ .

**Claim 2:** For  $k = 0, \dots, k_0$ ,  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \eta$ .

To prove the claim, we show by induction that  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \varepsilon\kappa^k$  (which is bounded above by  $\eta$  provided that  $k \leq k_0$ ). For  $k = 0$ , by assumption  $\|\mathbf{w}^{(0)} - \mathbf{w}^*\| \leq \varepsilon = \varepsilon\kappa^0$ , as required. For the inductive step, suppose that  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \varepsilon\kappa^k$  for  $k < k_0$ . Now, using  $\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \|G(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\|$  and the fact that  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \eta$ ,

$$\|\mathbf{w}^{(k+1)} - \mathbf{w}^*\| \leq \|G(\mathbf{w}^{(k)})\| \|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \kappa(\varepsilon\kappa^k) = \varepsilon\kappa^{k+1},$$

and the result now follows by induction.

**Claim 3:** For  $k = 0, \dots, k_0$ ,  $\mu_{\bar{A}}^{(k)}$  is monotonically nonincreasing in  $k$ , and  $\mu_{\bar{A}}^{(k_0)} = \mathbf{0}$  (which is equal to  $\mu_{\bar{A}}^*$ ).

By claim 2,  $\mathbf{g}_{\bar{A}}(\mathbf{x}^{(k)}) \leq -\delta\mathbf{e}$  for all  $k = 0, \dots, k_0$ . Hence, for  $k < k_0$ ,

$$\begin{aligned}\mu_{\bar{A}}^{(k+1)} &= [\mu_{\bar{A}}^{(k)} + \alpha\mathbf{g}_{\bar{A}}(\mathbf{x}^{(k)})]_+ \\ &\leq [\mu_{\bar{A}}^{(k)} - \alpha\delta\mathbf{e}]_+ \\ &\leq \mu_{\bar{A}}^{(k)},\end{aligned}$$

which establishes nonincreasing monotonicity.

To show that  $\mu_{\bar{A}}^{(k_0)} = \mathbf{0}$ , suppose that for some nonactive component  $l$ , we have  $\mu_l^{(k_0)} > 0$ . By the monotonicity above,  $\mu_l^{(k)} > 0$  for  $k = 0, \dots, k_0$ . Hence,

$$\begin{aligned}\mu_l^{(k_0)} &= \mu_l^{(k_0-1)} + \alpha g_l(\mathbf{x}^{(k_0-1)}) \\ &= \mu_l^{(0)} + \sum_{k=0}^{k_0-1} \alpha g_l(\mathbf{x}^{(k)}).\end{aligned}$$

But by claim 2,  $g_l(\mathbf{x}^{(k)}) \leq -\delta$  for all  $k = 0, \dots, k_0 - 1$ . Hence,  $\mu_l^{(k_0)} \leq \varepsilon - k_0\alpha\delta \leq 0$ , which is a contradiction.

Finally, we will state and prove claim 4, which completes the proof of the theorem.

**Claim 4:** For  $k \geq k_0$ , we have  $\mu_{\bar{A}}^{(k)} = \mathbf{0} = \mu_{\bar{A}}^*$ ,  $\|\mathbf{w}_A^{(k+1)} - \mathbf{w}_A^*\| \leq \kappa_A \|\mathbf{w}_A^{(k)} - \mathbf{w}_A^*\|$ , and  $\|\mathbf{w}^{(k)} - \mathbf{w}^*\| \leq \eta$ .

We use induction. For  $k = k_0$ , we have  $\|\mathbf{w}(k_0) - \mathbf{w}^*\| \leq \eta$  by claim 2,  $\mu_{\bar{A}}^{(k_0)} = \mathbf{0}$  by claim 3. Hence,

$$\mathbf{w}_A^{(k_0+1)} = \mathbf{\Pi}[\mathbf{U}_A(\mathbf{w}_A^{(k_0)}) + \alpha D\mathbf{g}_{\bar{A}}(\mathbf{x}^{(k_0)})^\top \mu_{\bar{A}}^{(k_0)}] = \mathbf{\Pi}[\mathbf{U}_A(\mathbf{w}_A^{(k_0)})].$$

Because  $\mu_A^* = \mathbf{0}$ , it is, similarly, also true that  $w_A^* = \Pi[U_A(w_A^*)]$ . Thus,

$$\begin{aligned}\|w_A^{(k_0+1)} - w_A^*\| &= \|\Pi[U_A(w_A^{(k_0)})] - \Pi[U_A(w_A^*)]\| \\ &\leq \|U_A(w_A^{(k_0)}) - U_A(w_A^*)\| \\ &\leq \|G_A(w_A^{(k_0)})\| \|w_A^{(k_0)} - w_A^*\|,\end{aligned}$$

where  $\|G_A(w_A^{(k_0)})\| \leq \kappa_A$  because  $\|w^{(k_0)} - w^*\| \leq \eta$ . Hence,  $\|w_A^{(k_0+1)} - w_A^*\| \leq \kappa_A \|w_A^{(k_0)} - w_A^*\|$ , as required.

For the inductive step, suppose that the result holds for  $k \geq k_0$ . Now,  $g_{\bar{A}}(\mathbf{x}^{(k)}) \leq -\delta e$  and

$$\mu_{\bar{A}}^{(k+1)} = [\mu_{\bar{A}}^{(k)} + \alpha g_{\bar{A}}(\mathbf{x}^{(k)})]_+ \leq [\mathbf{0} - \alpha \delta e]_+ = \mathbf{0},$$

which implies that  $\mu_{\bar{A}}^{(k+1)} = \mathbf{0}$ . It follows that

$$\begin{aligned}w_A^{(k+2)} &= \Pi[U_A(w_A^{(k+1)}) + \alpha D g_{\bar{A}}(\mathbf{x}^{(k+1)})^\top \mu_{\bar{A}}^{(k+1)}] \\ &= \Pi[U_A(w_A^{(k+1)})],\end{aligned}$$

and now using the same argument as in the case of  $k = k_0$  above we get  $\|w_A^{(k+2)} - w_A^*\| \leq \kappa_A \|w_A^{(k+1)} - w_A^*\|$ . Finally,

$$\|w^{(k+1)} - w^*\| = \|w_A^{(k+1)} - w_A^*\| \leq \kappa_A \|w_A^{(k)} - w_A^*\| \leq \eta.$$

Because  $\kappa_A < 1$ , claim 4 implies that  $w^{(k)}$  converges to  $w^*$ , with at least a linear order of convergence. ■

An application of Lagrangian algorithms to a problem in decentralized rate control for sensor networks appears in [24], [25], and [93]. The proof above is based on [25].

### 23.5 Penalty Methods

Consider a general constrained optimization problem

$$\begin{aligned}&\text{minimize } f(\mathbf{x}) \\ &\text{subject to } \mathbf{x} \in \Omega\end{aligned}$$

We now discuss a method for solving this problem using techniques from unconstrained optimization. Specifically, we approximate the constrained optimization problem above by the unconstrained optimization problem

$$\text{minimize } f(\mathbf{x}) + \gamma P(\mathbf{x}),$$

where  $\gamma \in \mathbb{R}$  is a positive constant and  $P : \mathbb{R}^n \rightarrow \mathbb{R}$  is a given function. We then solve the associated unconstrained optimization problem and use the solution as an approximation to the minimizer of the original problem. The constant  $\gamma$  is called the *penalty parameter*, and the function  $P$  is called the *penalty function*. Formally, we define a penalty function as follows.

**Definition 23.1** A function  $P : \mathbb{R}^n \rightarrow \mathbb{R}$  is called a *penalty function* for the constrained optimization problem above if it satisfies the following three conditions:

1.  $P$  is continuous.
2.  $P(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathbb{R}^n$ .
3.  $P(\mathbf{x}) = 0$  if and only if  $\mathbf{x}$  is feasible (i.e.,  $\mathbf{x} \in \Omega$ ). ■

Clearly, for the unconstrained problem above to be a good approximation to the original problem, the penalty function  $P$  must be chosen appropriately. The role of the penalty function is to “penalize” points that are outside the feasible set.

To illustrate how we choose penalty functions, consider a constrained optimization problem of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, p, \end{aligned}$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i = 1, \dots, p$ . Considering only inequality constraints is not restrictive, because an equality constraint of the form  $\mathbf{h}(\mathbf{x}) = \mathbf{0}$  is equivalent to the inequality constraint  $\|\mathbf{h}(\mathbf{x})\|^2 \leq 0$  (however, see Exercise 21.25 for a caveat). For the constrained problem above, it is natural that the penalty function be defined in terms of the constraint functions  $g_1, \dots, g_p$ . A possible choice for  $P$  is

$$P(\mathbf{x}) = \sum_{i=1}^p g_i^+(\mathbf{x}),$$

where

$$g_i^+(\mathbf{x}) = \max\{0, g_i(\mathbf{x})\} = \begin{cases} 0 & \text{if } g_i(\mathbf{x}) \leq 0 \\ g_i(\mathbf{x}) & \text{if } g_i(\mathbf{x}) > 0. \end{cases}$$

We refer to this penalty function as the *absolute value penalty function*, because it is equal to  $\sum |g_i(\mathbf{x})|$ , where the summation is taken over all constraints that are violated at  $\mathbf{x}$ . We illustrate this penalty function in the following example.

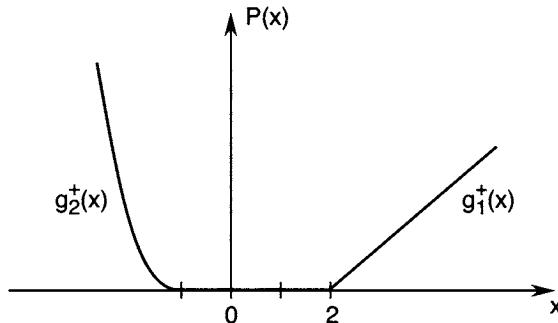


Figure 23.1  $g^+$  for Example 23.2.

**Example 23.2** Let  $g_1, g_2 : \mathbb{R} \rightarrow \mathbb{R}$  be defined by  $g_1(x) = x - 2$ ,  $g_2(x) = -(x + 1)^3$ . The feasible set defined by  $\{x \in \mathbb{R} : g_1(x) \leq 0, g_2(x) \leq 0\}$  is simply the interval  $[-1, 2]$ . In this example, we have

$$g_1^+(x) = \max\{0, g_1(x)\} = \begin{cases} 0 & \text{if } x \leq 2 \\ x - 2 & \text{otherwise,} \end{cases}$$

$$g_2^+(x) = \max\{0, g_2(x)\} = \begin{cases} 0 & \text{if } x \geq -1 \\ -(x + 1)^3 & \text{otherwise,} \end{cases}$$

and

$$P(x) = g_1^+(x) + g_2^+(x) = \begin{cases} x - 2 & \text{if } x > 2 \\ 0 & \text{if } -1 \leq x \leq 2 \\ -(x + 1)^3 & \text{if } x < -1. \end{cases}$$

Figure 23.1 provides a graphical illustration of  $g^+$  for this example. ■

The absolute value penalty function may not be differentiable at points  $\mathbf{x}$  where  $g_i(\mathbf{x}) = 0$ , as is the case at the point  $x = 2$  in Example 23.2 (notice, though, that in Example 23.2,  $P$  is differentiable at  $x = -1$ ). Therefore, in such cases we cannot use techniques for optimization that involve derivatives. A form of the penalty function that is guaranteed to be differentiable is the *Courant-Beltrami penalty function*, given by

$$P(\mathbf{x}) = \sum_{i=1}^p (g_i^+(\mathbf{x}))^2.$$

In the following discussion we do not assume any particular form of the penalty function  $P$ . We only assume that  $P$  satisfies conditions 1 to 3 given in Definition 23.1.

The penalty function method for solving constrained optimization problems involves constructing and solving an associated unconstrained optimization

problem and using the solution to the unconstrained problem as the solution to the original constrained problem. Of course, the solution to the unconstrained problem (the approximated solution) may not be exactly equal to the solution to the constrained problem (the true solution). Whether or not the solution to the unconstrained problem is a good approximation to the true solution depends on the penalty parameter  $\gamma$  and the penalty function  $P$ . We would expect that the larger the value of the penalty parameter  $\gamma$ , the closer the approximated solution will be to the true solution, because points that violate the constraints are penalized more heavily. Ideally, in the limit as  $\gamma \rightarrow \infty$ , the penalty method should yield the true solution to the constrained problem. In the remainder of this section, we analyze this property of the penalty function method.

**Example 23.3** Consider the problem

$$\begin{aligned} & \text{minimize } \mathbf{x}^\top \mathbf{Q} \mathbf{x} \\ & \text{subject to } \|\mathbf{x}\|^2 = 1, \end{aligned}$$

where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ .

- a. Using the penalty function  $P(\mathbf{x}) = (\|\mathbf{x}\|^2 - 1)^2$  and penalty parameter  $\gamma$ , write down an unconstrained optimization problem whose solution  $\mathbf{x}_\gamma$  approximates the solution to this problem.
- b. Show that for any  $\gamma$ ,  $\mathbf{x}_\gamma$  is an eigenvector of  $\mathbf{Q}$ .
- c. Show that  $\|\mathbf{x}_\gamma\|^2 - 1 = O(1/\gamma)$  as  $\gamma \rightarrow \infty$ .

*Solution:*

- a. The unconstrained problem based on the given penalty function is

$$\text{minimize } \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \gamma(\|\mathbf{x}\|^2 - 1)^2.$$

- b. By the FONC,  $\mathbf{x}_\gamma$  satisfies

$$2\mathbf{Q}\mathbf{x}_\gamma + 4\gamma(\|\mathbf{x}_\gamma\|^2 - 1)\mathbf{x}_\gamma = 0.$$

Rearranging, we obtain

$$\mathbf{Q}\mathbf{x}_\gamma = 2\gamma(1 - \|\mathbf{x}_\gamma\|^2)\mathbf{x}_\gamma = \lambda_\gamma \mathbf{x}_\gamma,$$

where  $\lambda_\gamma$  is a scalar. Hence,  $\mathbf{x}_\gamma$  is an eigenvector of  $\mathbf{Q}$ . (This agrees with Example 20.8.)

- c. Now,  $\lambda_\gamma = 2\gamma(1 - \|\mathbf{x}_\gamma\|^2) \leq \lambda_{\max}$ , where  $\lambda_{\max}$  is the largest eigenvalue of  $\mathbf{Q}$ . Hence,  $\|\mathbf{x}_\gamma\|^2 - 1 = -\lambda_{\max}/(2\gamma) = O(1/\gamma)$  as  $\gamma \rightarrow \infty$ . ■

We now analyze the penalty method in a more general setting. In our analysis, we adopt the following notation. Denote by  $\mathbf{x}^*$  a solution (global minimizer) to the problem. Let  $P$  be a penalty function for the problem. For each  $k = 1, 2, \dots$ , let  $\gamma_k \in \mathbb{R}$  be a given positive constant. Define an associated function  $q(\gamma_k, \cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$  by

$$q(\gamma_k, \mathbf{x}) = f(\mathbf{x}) + \gamma_k P(\mathbf{x}).$$

For each  $k$ , we can write the following associated unconstrained optimization problem:

$$\text{minimize } q(\gamma_k, \mathbf{x}).$$

Denote by  $\mathbf{x}^{(k)}$  a minimizer of  $q(\gamma_k, \mathbf{x})$ . The following technical lemma describes certain useful relationships between the constrained problem and the associated unconstrained problems.

**Lemma 23.4** *Suppose that  $\{\gamma_k\}$  is a nondecreasing sequence; that is, for each  $k$ , we have  $\gamma_k \leq \gamma_{k+1}$ . Then, for each  $k$  we have*

1.  $q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) \geq q(\gamma_k, \mathbf{x}^{(k)})$ .
2.  $P(\mathbf{x}^{(k+1)}) \leq P(\mathbf{x}^{(k)})$ .
3.  $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$ .
4.  $f(\mathbf{x}^*) \geq q(\gamma_k, \mathbf{x}^{(k)}) \geq f(\mathbf{x}^{(k)})$ .

□

*Proof.* We first prove part 1. From the definition of  $q$  and the fact that  $\{\gamma_k\}$  is an increasing sequence, we have

$$q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k+1)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}).$$

Now, because  $\mathbf{x}^{(k)}$  is a minimizer of  $q(\gamma_k, \mathbf{x})$ ,

$$q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}).$$

Combining the above, we get part 1.

We next prove part 2. Because  $\mathbf{x}^{(k)}$  and  $\mathbf{x}^{(k+1)}$  minimize  $q(\gamma_k, \mathbf{x})$  and  $q(\gamma_{k+1}, \mathbf{x})$ , respectively, we can write

$$q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}),$$

$$q(\gamma_{k+1}, \mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k+1)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) + \gamma_{k+1} P(\mathbf{x}^{(k)}).$$

Adding the inequalities above yields

$$\gamma_k P(\mathbf{x}^{(k)}) + \gamma_{k+1} P(\mathbf{x}^{(k+1)}) \leq \gamma_{k+1} P(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k+1)}).$$

Rearranging, we get

$$(\gamma_{k+1} - \gamma_k) P(\mathbf{x}^{(k+1)}) \leq (\gamma_{k+1} - \gamma_k) P(\mathbf{x}^{(k)}).$$

We know by assumption that  $\gamma_{k+1} \geq \gamma_k$ . If  $\gamma_{k+1} > \gamma_k$ , then we get  $P(\mathbf{x}^{(k+1)}) \leq P(\mathbf{x}^{(k)})$ . If, on the other hand,  $\gamma_{k+1} = \gamma_k$ , then clearly  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$  and so  $P(\mathbf{x}^{(k+1)}) = P(\mathbf{x}^{(k)})$ . Therefore, in either case, we arrive at part 2.

We now prove part 3. Because  $\mathbf{x}^{(k)}$  is a minimizer of  $q(\gamma_k, \mathbf{x})$ , we obtain

$$q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)}) + \gamma_k P(\mathbf{x}^{(k+1)}).$$

Therefore,

$$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)}) + \gamma_k (P(\mathbf{x}^{(k)}) - P(\mathbf{x}^{(k+1)})).$$

From part 2 we have  $P(\mathbf{x}^{(k)}) - P(\mathbf{x}^{(k+1)}) \geq 0$ , and  $\gamma_k > 0$  by assumption; therefore, we get

$$f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)}).$$

Finally, we now prove part 4. Because  $\mathbf{x}^{(k)}$  is a minimizer of  $q(\gamma_k, \mathbf{x})$ , we get

$$f(\mathbf{x}^*) + \gamma_k P(\mathbf{x}^*) \geq q(\gamma_k, \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}).$$

Because  $\mathbf{x}^*$  is a minimizer for the constrained optimization problem, we have  $P(\mathbf{x}^*) = 0$ . Therefore,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}^{(k)}) + \gamma_k P(\mathbf{x}^{(k)}).$$

Because  $P(\mathbf{x}^{(k)}) \geq 0$  and  $\gamma_k \geq 0$ ,

$$f(\mathbf{x}^*) \geq q(\gamma_k, \mathbf{x}^{(k)}) \geq f(\mathbf{x}^{(k)}),$$

which completes the proof. ■

With the above lemma, we are now ready to prove the following theorem.

**Theorem 23.4** Suppose that the objective function  $f$  is continuous and  $\gamma_k \rightarrow \infty$  as  $k \rightarrow \infty$ . Then, the limit of any convergent subsequence of the sequence  $\{\mathbf{x}^{(k)}\}$  is a solution to the constrained optimization problem. □

*Proof.* Suppose that  $\{\mathbf{x}^{(m_k)}\}$  is a convergent subsequence of the sequence  $\{\mathbf{x}^{(k)}\}$ . (See Section 5.1 for a discussion of sequences and subsequences.) Let  $\hat{\mathbf{x}}$  be the limit of  $\{\mathbf{x}^{(m_k)}\}$ . By Lemma 23.4, the sequence  $\{q(\gamma_k, \mathbf{x}^{(k)})\}$  is non-decreasing and bounded above by  $f(\mathbf{x}^*)$ . Therefore, the sequence  $\{q(\gamma_k, \mathbf{x}^{(k)})\}$  has a limit  $q^* = \lim_{k \rightarrow \infty} q(\gamma_k, \mathbf{x}^{(k)})$  such that  $q^* \leq f(\mathbf{x}^*)$  (see Theorem 5.3). Because the function  $f$  is continuous and  $f(\mathbf{x}^{(m_k)}) \leq f(\mathbf{x}^*)$  by Lemma 23.4, we have

$$\lim_{k \rightarrow \infty} f(\mathbf{x}^{(m_k)}) = f\left(\lim_{k \rightarrow \infty} \mathbf{x}^{(m_k)}\right) = f(\hat{\mathbf{x}}) \leq f(\mathbf{x}^*).$$

Because the sequences  $\{f(\mathbf{x}^{(m_k)})\}$  and  $\{q(\gamma_{m_k}, \mathbf{x}^{(m_k)})\}$  both converge, the sequence  $\{\gamma_{m_k} P(\mathbf{x}^{(m_k)})\} = \{q(\gamma_{m_k}, \mathbf{x}^{(m_k)}) - f(\mathbf{x}^{(m_k)})\}$  also converges, with

$$\lim_{k \rightarrow \infty} \gamma_{m_k} P(\mathbf{x}^{(m_k)}) = q^* - f(\hat{\mathbf{x}}).$$

By Lemma 23.4, the sequence  $\{P(\mathbf{x}^{(k)})\}$  is nonincreasing and bounded from below by 0. Therefore,  $\{P(\mathbf{x}^{(k)})\}$  converges (again see Theorem 5.3), and hence so does  $\{P(\mathbf{x}^{(m_k)})\}$ . Because  $\gamma_{m_k} \rightarrow \infty$  we conclude that

$$\lim_{k \rightarrow \infty} P(\mathbf{x}^{(m_k)}) = 0.$$

By continuity of  $P$ , we have

$$0 = \lim_{k \rightarrow \infty} P(\mathbf{x}^{(m_k)}) = P\left(\lim_{k \rightarrow \infty} \mathbf{x}^{(m_k)}\right) = P(\hat{\mathbf{x}}),$$

and hence  $\hat{\mathbf{x}}$  is a feasible point. Because  $f(\mathbf{x}^*) \geq f(\hat{\mathbf{x}})$  from above, we conclude that  $\hat{\mathbf{x}}$  must be a solution to the constrained optimization problem. ■

If we perform an infinite number of minimization runs, with the penalty parameter  $\gamma_k \rightarrow \infty$ , then Theorem 23.4 ensures that the limit of any convergent subsequence is a minimizer  $\mathbf{x}^*$  to the original constrained optimization problem. There is clearly a practical limitation in applying this theorem. It is certainly desirable to find a minimizer to the original constrained optimization problem using a *single* minimization run for the unconstrained problem that approximates the original problem using a penalty function. In other words, we desire an exact solution to the original constrained problem by solving the associated unconstrained problem [minimize  $f(\mathbf{x}) + \gamma P(\mathbf{x})$ ] with a finite  $\gamma > 0$ . It turns out that indeed this can be accomplished, in which case we say that the penalty function is *exact*. However, it is necessary that exact penalty functions be nondifferentiable, as shown in [10], and illustrated in the following example.

**Example 23.4** Consider the problem

$$\begin{aligned} &\text{minimize} && f(x) \\ &\text{subject to} && x \in [0, 1], \end{aligned}$$

where  $f(x) = 5 - 3x$ . Clearly, the solution is  $x^* = 1$ .

Suppose that we use the penalty method to solve the problem, with a penalty function  $P$  that is differentiable at  $x^* = 1$ . Then,  $P'(x^*) = 0$ , because  $P(x) = 0$  for all  $x \in [0, 1]$ . Hence, if we let  $g = f + \gamma P$ , then  $g'(x^*) = f'(x^*) + \gamma P'(x^*) \neq 0$  for all finite  $\gamma > 0$ . Hence,  $x^* = 1$  does not satisfy the first-order necessary condition to be a local minimizer of  $g$ . Thus,  $P$  is not an exact penalty function. ■

Here, we prove a result on the necessity of nondifferentiability of exact penalty functions for a special class of problems.

**Proposition 23.4** *Consider the problem*

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

with  $\Omega \subset \mathbb{R}^n$  convex. Suppose that the minimizer  $\mathbf{x}^*$  lies on the boundary of  $\Omega$  and there exists a feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$  such that  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) > 0$ . If  $P$  is an exact penalty function, then  $P$  is not differentiable at  $\mathbf{x}^*$ .  $\square$

*Proof.* We use contraposition. Suppose that  $P$  is differentiable at  $\mathbf{x}^*$ . Then,  $\mathbf{d}^\top \nabla P(\mathbf{x}^*) = 0$ , because  $P(\mathbf{x}) = 0$  for all  $\mathbf{x} \in \Omega$ . Hence, if we let  $g = f + \gamma P$ , then  $\mathbf{d}^\top \nabla g(\mathbf{x}^*) > 0$  for all finite  $\gamma > 0$ , which implies that  $\nabla g(\mathbf{x}^*) \neq \mathbf{0}$ . Hence,  $\mathbf{x}^*$  is not a local minimizer of  $g$ , and thus  $P$  is not an exact penalty function.  $\blacksquare$

Note that the result of Proposition 23.4 does not hold if we remove the assumption that  $\mathbf{d}^\top \nabla f(\mathbf{x}^*) > 0$ . Indeed, consider a convex problem where  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Choose  $P$  to be differentiable. Clearly, in this case we have  $\nabla g(\mathbf{x}^*) = \nabla f(\mathbf{x}^*) + \gamma \nabla P(\mathbf{x}^*) = \mathbf{0}$ . The function  $P$  is therefore an exact penalty function, although differentiable.

For further reading on the subject of optimization of nondifferentiable functions, see, for example, [38]. References [11] and [96] provide further discussions on the penalty method, including nondifferentiable exact penalty functions. These references also discuss exact penalty methods involving differentiable functions; these methods go beyond the elementary type of penalty method introduced in this chapter.

## EXERCISES

**23.1** Consider the constrained optimization problem

$$\begin{aligned} &\text{maximize} && f(\mathbf{x}) \\ &\text{subject to} && \|\mathbf{x}\| = 1, \end{aligned}$$

where  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x}$  and  $\mathbf{Q} = \mathbf{Q}^\top$ . We wish to apply a fixed-step-size projected gradient algorithm to this problem:

$$\mathbf{x}^{(k+1)} = \mathbf{\Pi}[\mathbf{x}^{(k)} + \alpha \nabla f(\mathbf{x}^{(k)})],$$

where  $\alpha > 0$  and  $\mathbf{\Pi}$  is the usual projection operator defined by  $\mathbf{\Pi}[\mathbf{x}] = \arg \min_{\mathbf{z} \in \Omega} \|\mathbf{z} - \mathbf{x}\|$  and  $\Omega$  is the constraint set.

- a. Find a simple formula for  $\mathbf{\Pi}[\mathbf{x}]$  in this problem (an explicit expression in terms of  $\mathbf{x}$ ), assuming that  $\mathbf{x} \neq \mathbf{0}$ .
- b. For the remainder of the question, suppose that

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}.$$

Find the solution(s) to this optimization problem.

- c. Let  $y^{(k)} = x_1^{(k)}/x_2^{(k)}$ . Derive an expression for  $y^{(k+1)}$  in terms of  $y^{(k)}$  and  $\alpha$ .
- d. Assuming that  $x_2^{(0)} \neq 0$ , use parts b and c to show that for any  $\alpha > 0$ ,  $\mathbf{x}^{(k)}$  converges to a solution to the optimization problem (i.e., the algorithm works).
- e. In part d, what if  $x_2^{(0)} = 0$ ?

### 23.2 Consider the problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$  and  $\mathbf{c} \in \mathbb{R}^n$  is a given nonzero vector. (Linear programming is a special case of this problem.) We wish to apply a fixed-step-size projected gradient algorithm

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \nabla f(\mathbf{x}^{(k)})],$$

where, as usual,  $\Pi$  is the projection operator onto  $\Omega$  (assume that for any  $\mathbf{y}$ ,  $\Pi[\mathbf{y}] = \arg \min_{\mathbf{x} \in \Omega} \|\mathbf{y} - \mathbf{x}\|^2$  is unique).

- a. Suppose that for some  $k$ ,  $\mathbf{x}^{(k)}$  is a global minimizer of the problem. Is it necessarily the case that  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ ? Explain fully.
- b. Suppose that for some  $k$ ,  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$ . Is it necessarily the case that  $\mathbf{x}^{(k)}$  is a local minimizer of the problem? Explain fully.

### 23.3 Consider the optimization problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f \in C^1$ , and  $\Omega = [-1, 1]^2 = \{\mathbf{x} : -1 \leq x_i \leq 1, i = 1, 2\}$ . Consider the projected steepest descent algorithm applied to this problem:

$$\mathbf{x}^{(k+1)} = \Pi[\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})],$$

where  $\Pi$  represents the projection operator with respect to  $\Omega$  and  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$ . Our goal is to prove the following statement:

$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)}$  if and only if  $\mathbf{x}^{(k)}$  satisfies the first-order necessary condition.

We will do this in two parts.

- a. Prove the statement above for the case where  $\mathbf{x}^{(k)}$  is an interior point of  $\Omega$ .

- b. Prove the statement for the case where  $\mathbf{x}^{(k)}$  is a boundary point of  $\Omega$ .

*Hint:* Consider two further subcases: (i)  $\mathbf{x}^{(k)}$  is a corner point, and (ii)  $\mathbf{x}^{(k)}$  is not a corner point. For subcase (i) it suffices to take  $\mathbf{x}^{(k)} = [1, 1]^\top$ . For subcase (ii) it suffices to take  $\mathbf{x}^{(k)} \in \{\mathbf{x} : x_1 = 1, -1 < x_2 < 1\}$ .

**23.4** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ ,  $\text{rank } \mathbf{A} = m$ , and  $\mathbf{b} \in \mathbb{R}^m$ . Define  $\Omega = \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$  and let  $\mathbf{x}_0 \in \Omega$ . Show that for any  $\mathbf{y} \in \mathbb{R}^n$ ,

$$\Pi[\mathbf{x}_0 + \mathbf{y}] = \mathbf{x}_0 + \mathbf{P}\mathbf{y},$$

where  $\mathbf{P} = \mathbf{I} - \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{A}$ .

*Hint:* Use Exercise 6.7 and Example 12.5.

**23.5** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be given by  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{x}^\top \mathbf{c}$ , where  $\mathbf{Q} = \mathbf{Q}^\top > 0$ . We wish to minimize  $f$  over  $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ , where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ , and  $\text{rank } \mathbf{A} = m$ . Show that the projected steepest descent algorithm for this case takes the form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\mathbf{g}^{(k)\top} \mathbf{P} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)\top} \mathbf{P} \mathbf{Q} \mathbf{P} \mathbf{g}^{(k)}} \mathbf{P} \mathbf{g}^{(k)},$$

where

$$\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{c},$$

and  $\mathbf{P} = \mathbf{I}_n - \mathbf{A}^\top(\mathbf{AA}^\top)^{-1}\mathbf{A}$ .

**23.6** Consider the problem

$$\begin{aligned} &\text{minimize} && \frac{1}{2} \|\mathbf{x}\|^2 \\ &\text{subject to} && \mathbf{Ax} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $m < n$ , and  $\text{rank } \mathbf{A} = m$ . Show that if  $\mathbf{x}^{(0)} \in \{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ , then the projected steepest descent algorithm converges to the solution in one step.

**23.7** Show that in the projected steepest descent algorithm, we have that for each  $k$ :

- a.  $\mathbf{g}^{(k+1)\top} \mathbf{P} \mathbf{g}^{(k)} = 0$ .

- b. The vector  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$  is orthogonal to the vector  $\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}$ .

**23.8** Consider the optimization problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega, \end{aligned}$$

where  $\Omega \subset \mathbb{R}^n$ . Suppose that we apply the *penalty method* to this problem, which involves solving an associated unconstrained optimization problem with penalty function  $P$  and penalty parameter  $\gamma > 0$ .

- a. Write down the unconstrained problem associated with penalty function  $P$  and penalty parameter  $\gamma$ .
- b. Let  $\mathbf{x}^*$  be a global minimizer of the given constrained problem, and let  $\mathbf{x}^\gamma$  be a global minimizer of the associated unconstrained optimization problem (in part a) with penalty parameter  $\gamma$ . Show that if  $\mathbf{x}^\gamma \notin \Omega$ , then  $f(\mathbf{x}^\gamma) < f(\mathbf{x}^*)$ .

**23.9** Use the penalty method to solve the following problem:

$$\begin{aligned} & \text{minimize} && x_1^2 + 2x_2^2 \\ & \text{subject to} && x_1 + x_2 = 3. \end{aligned}$$

*Hint:* Use the penalty function  $P(x) = (x_1 + x_2 - 3)^2$ . The solution you find must be exact, not approximate.

**23.10** Consider the simple optimization problem

$$\begin{aligned} & \text{minimize} && x \\ & \text{subject to} && x \geq a, \end{aligned}$$

where  $a \in \mathbb{R}$ . Suppose that we use the penalty method to solve this problem, with penalty function

$$P(x) = (\max\{a - x, 0\})^2$$

(the *Courant-Beltrami penalty function*). Given a number  $\varepsilon > 0$ , find the smallest value of the penalty parameter  $\gamma$  such that the solution obtained using the penalty method is no further than  $\varepsilon$  from the true solution to the given problem. (Think of  $\varepsilon$  as the desired accuracy.)

**23.11** Consider the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{x}\|^2 \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \leq n$ , and  $\text{rank } \mathbf{A} = m$ . Let  $\mathbf{x}^*$  be the solution. Suppose that we solve the problem using the penalty method, with the penalty function

$$P(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|^2.$$

Let  $\mathbf{x}_\gamma^*$  be the solution to the associated unconstrained problem with the penalty parameter  $\gamma > 0$ ; that is,  $\mathbf{x}_\gamma^*$  is the solution to

$$\text{minimize } \frac{1}{2} \|\mathbf{x}\|^2 + \gamma \|\mathbf{Ax} - \mathbf{b}\|^2.$$

- a. Suppose that

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \end{bmatrix}, \quad \mathbf{b} = [1].$$

Verify that  $\mathbf{x}_\gamma^*$  converges to the solution  $\mathbf{x}^*$  of the original constrained problem as  $\gamma \rightarrow \infty$ .

- b. Prove that  $\mathbf{x}_\gamma^* \rightarrow \mathbf{x}^*$  as  $\gamma \rightarrow \infty$  holds in general.

*Hint:* Use the following result: There exist orthogonal matrices  $\mathbf{U} \in \mathbb{R}^{m \times m}$  and  $\mathbf{V}^\top \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{A} = \mathbf{U} [\mathbf{S}, \mathbf{O}] \mathbf{V}^\top,$$

where

$$\mathbf{S} = \text{diag} \left( \sqrt{\lambda_1(\mathbf{A}\mathbf{A}^\top)}, \dots, \sqrt{\lambda_m(\mathbf{A}\mathbf{A}^\top)} \right)$$

is a diagonal matrix with diagonal elements that are the square roots of the eigenvalues of  $\mathbf{A}\mathbf{A}^\top$ .

The result above is called the *singular value decomposition* (see, e.g., [62, p. 411]).



## CHAPTER 24

---

# MULTIOBJECTIVE OPTIMIZATION

---

### 24.1 Introduction

When an optimization problem involves only one objective function, it is a single-objective optimization. Most engineering problems require the designer to optimize a number of conflicting objectives. The objectives are in conflict with each other if an improvement in one objective leads to deterioration in another. Multiobjective problems in which there is competition between objectives may have no single, unique optimal solution. Multiobjective optimization problems are also referred to as multicriteria or vector optimization problems. We can formulate a multiobjective optimization problem as follows: Find a decision variable that satisfies the given constraints and optimizes a vector function whose components are objective functions. Formally, the mul-

tiobjective optimization problem is stated as follows:

$$\begin{aligned} \text{minimize } \mathbf{f}(\mathbf{x}) = & \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_\ell(x_1, x_2, \dots, x_n) \end{bmatrix} \\ \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$  and  $\Omega \subset \mathbb{R}^n$ . For example, the constraint set  $\Omega$  can have the form

$$\Omega = \{\mathbf{x} : \mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}\},$$

where

$$\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p, \quad m \leq n.$$

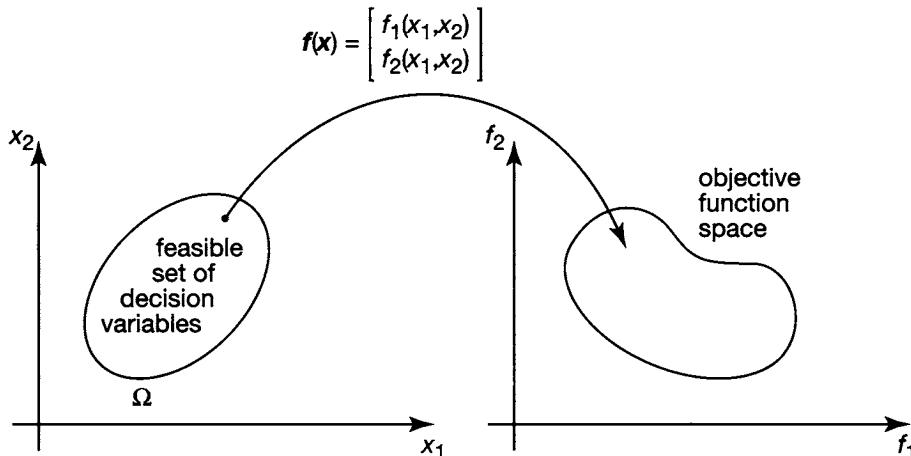
In general, we may have three different types of multiobjective optimization problems:

- Minimize all the objective functions.
- Maximize all the objective functions.
- Minimize some and maximize others.

However, as usual, any of these can be converted into an equivalent minimization problem.

## 24.2 Pareto Solutions

Note that multiobjective function assigns to each decision variable a multi-objective vector value in the objective function space. A graphical illustration of this statement is illustrated in Figures 24.1 and 24.2. In Figure 24.1 the decision variable is a point  $\mathbf{x} \in \mathbb{R}^2$  while the vector of objective functions is given by  $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ . In Figure 24.2 the decision variable is a point  $\mathbf{x} \in \mathbb{R}^2$  while the vector of objective functions is  $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ . In single-objective optimization problems our goal is to find a single solution, where we focus mainly on the decision variable space. On the other hand, in multiobjective problems we are usually more interested in the objective space. As pointed out by Miettinen [92, p. 11], multiobjective problems are in a sense ill-defined because there is no natural ordering in the objective space. Miettinen [92] illustrates this statement with the following simple example. One can say that  $[1, 1]^\top$  is less than  $[3, 3]^\top$ . But how do we compare  $[1, 3]^\top$  and  $[3, 1]^\top$ ? In general, in multiobjective optimization problems our goal is to find good compromises. Roughly speaking, in a multiobjective optimization problem, a solution is optimal if there exists no other solution, within the feasible set,



**Figure 24.1** Two-dimensional illustration of a multiobjective vector function assigning to each decision variable a multiobjective vector value.

that gives improved performance with regard to all the objectives. A formal definition of an optimal point for a multiobjective optimization problem was proposed by Francis Y. Edgeworth in 1881 and generalized by Vilfredo Pareto in 1896. It is customary now to refer to an optimal point of a multiobjective optimization problem as the Pareto minimizer, whose formal definition is given next.

**Definition 24.1** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$  and  $\mathbf{x} \in \Omega$  be given. For the optimization problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega \end{aligned}$$

a point  $\mathbf{x}^* \in \Omega$  is called a *Pareto minimizer* if there exists no  $\mathbf{x} \in \Omega$  such that for  $i = 1, 2, \dots, \ell$ ,

$$f_i(\mathbf{x}) \leq f_i(\mathbf{x}^*)$$

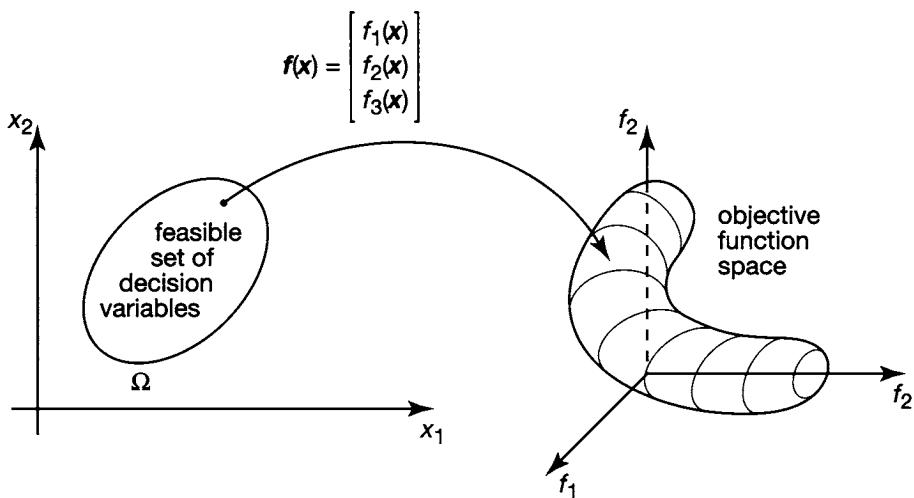
and for at least one  $i$ ,

$$f_i(\mathbf{x}) < f_i(\mathbf{x}^*)$$

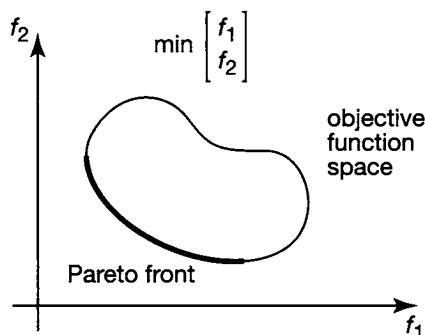
■

In other words, the point  $\mathbf{x}^*$  is a Pareto minimizer, or a nondominated solution, if there exists no other feasible decision variable  $\mathbf{x}$  that would decrease some objectives without causing simultaneous increase in at least one other variable.

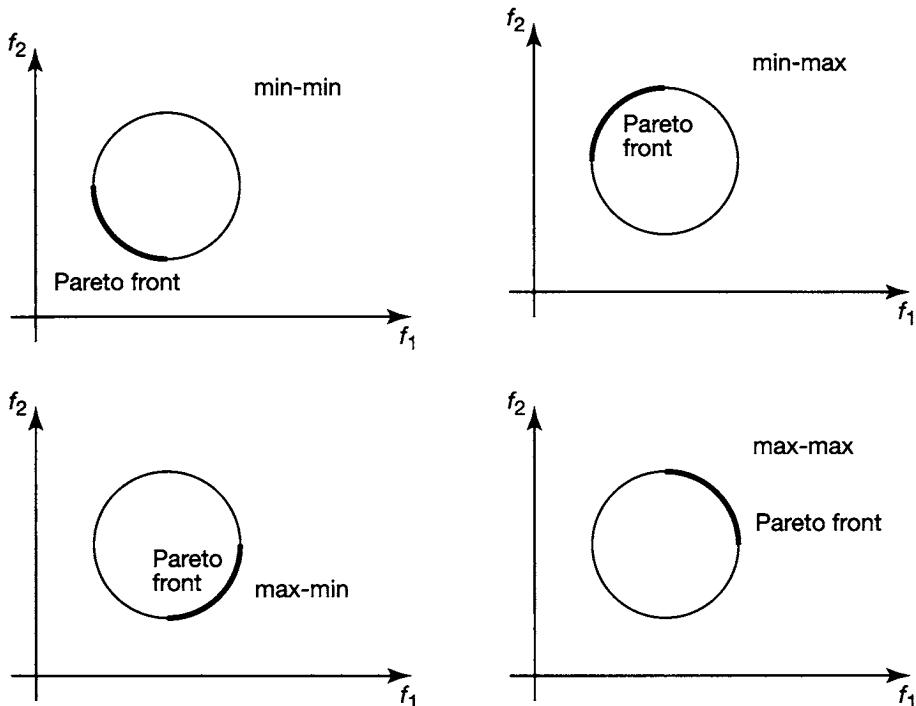
The set of Pareto minimizers (optimizers) is called the *Pareto front*, as illustrated in Figure 24.3. Most multiobjective optimization algorithms use



**Figure 24.2** Three-dimensional illustration of a multiobjective vector function assigning to each decision variable a multiobjective vector function value.



**Figure 24.3** The Pareto front is marked with a heavy line.



**Figure 24.4** Pareto fronts for four possible cases of two-objective optimization.

the concept of domination. A solution is said to be nondominated if it is Pareto optimal.

In Figure 24.4 we show different combinations of two-objective optimization and the corresponding Pareto fronts. In particular, in the upper left, we show the Pareto front for the case when we are minimizing both components of the objective function vector, which we represent by “min-min.” Similarly, “min-max” represents the case when we are minimizing the first objective function and maximizing the second; and so forth.

### 24.3 Computing the Pareto Front

When computing the Pareto front, two solutions are compared and the dominated solution is eliminated from the set of candidates of Pareto optimizers. Thus, the Pareto front consists of nondominated solutions.

To proceed, we need some notation. Let

$$\mathbf{x}^{*r} = [x_1^{*r}, x_2^{*r}, \dots, x_n^{*r}]^\top$$

be the  $r$ th candidate Pareto optimal solution,  $r = 1, 2, \dots, R$ , where  $R$  is the number of current candidate Pareto solutions. Let

$$\mathbf{f}(\mathbf{x}^{*r}) = [f_1(\mathbf{x}^{*r}), f_2(\mathbf{x}^{*r}), \dots, f_\ell(\mathbf{x}^{*r})]^\top$$

be the corresponding value of the objective function vector. For any new solution candidate  $\mathbf{x}^j$ , we evaluate the objective function vector  $\mathbf{f}(\mathbf{x}^j)$ . We then compare the new solution candidate with the existing Pareto solutions. We need to consider three cases:

- $\mathbf{x}^j$  dominates at least one candidate solution.
- $\mathbf{x}^j$  does not dominate any existing candidate solutions.
- $\mathbf{x}^j$  is dominated by a candidate solution.

If  $\mathbf{x}^j$  dominates at least one candidate solution, we delete the dominated solutions from the set and add the new solution  $\mathbf{x}^j$  to the set of candidates. In the second case, when the new candidate solution  $\mathbf{x}^j$  does not dominate any of the existing candidate Pareto solutions, add this new Pareto solution to the set of candidate Pareto solutions. Finally, in the third case, when the new candidate solution is dominated by at least one of the existing candidate Pareto solutions, we do not change the set of the existing candidate Pareto solutions.

**Example 24.1** Consider the two-objective minimization problem whose data are as follows:

$\mathbf{x}^{(i)\top}$	$\mathbf{f}(\mathbf{x}^{(i)})^\top$
[5, 6]	[30, 45]
[4, 5]	[22, 29]
[3, 7]	[19, 53]
[6, 8]	[41, 75]
[1, 4]	[13, 45]
[6, 7]	[42, 55]
[2, 5]	[37, 46]
[3, 6]	[28, 37]
[2, 7]	[12, 51]
[4, 7]	[41, 67]

Suppose that we wish to find nondominated pairs for this problem. Recall that a point  $\mathbf{x}^*$  is a nondominated point if for all  $i$  and all  $\mathbf{x}$ ,

$$f_i(\mathbf{x}^*) \leq f_i(\mathbf{x}),$$

and at least for one component  $j$  of the objective vector, we have

$$f_j(\mathbf{x}^*) < f_j(\mathbf{x}).$$

To find the Pareto front, we start with the first pair as a candidate Pareto optimal solution and then compare the other pairs against this first pair, replacing the first pair as necessary. We then continue with the other pairs, building up a set of candidate Pareto solutions and modifying this set when appropriate. The result of the search gives the following Pareto optimal set:

$\mathbf{x}^{(i)\top}$	$\mathbf{f}(\mathbf{x}^{(i)})^\top$
[4, 5]	[22, 29]
[1, 4]	[13, 45]
[2, 7]	[12, 51]

■

We now discuss an algorithm for generating the Pareto front that implements the foregoing ideas. This algorithm is a minor modification of the algorithm of Osyczka [98, pp. 100–101]. We use the following notation. Let  $J$  be the number of candidate solutions to be checked for optimality, while  $R$  is the number of current candidate Pareto solutions. Recall that  $\ell$  is the number of objective functions, the dimension of the objective function vector, and  $n$  is the dimension of the decision space, that is, the number of components of  $\mathbf{x}$ . The algorithm consists of eight steps.

### Algorithm for Generating a Pareto Front

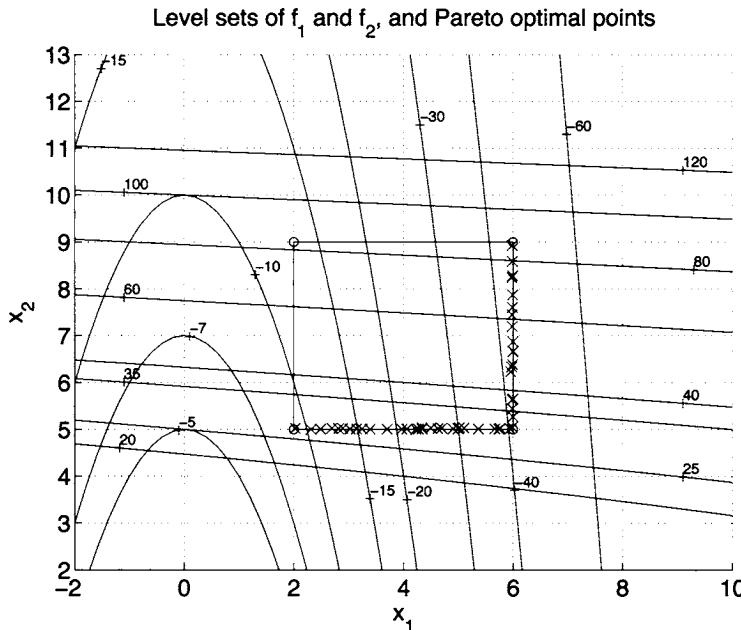
1. Generate an initial solution  $\mathbf{x}^1$  and evaluate  $\mathbf{f}^{*1} = \mathbf{f}(\mathbf{x}^1)$ . This first solution generated is taken as a candidate Pareto solution. Set initial indices  $R := 1$  and  $j := 1$ .
2. Set  $j := j + 1$ . If  $j \leq J$ , then generate solution  $\mathbf{x}^j$  and go to step 3. Otherwise, stop, because all the candidate solutions have already been considered.
3. Set  $r := 1$  and  $q := 0$  ( $q$  represents the number of eliminated solutions from the existing set of Pareto solutions).
4. If for all  $i = 1, 2, \dots, \ell$ ,
$$f_i(\mathbf{x}^j) < f_i(\mathbf{x}^{*r}),$$

then set  $q := q + 1$ ,  $\mathbf{f}^{*R} := \mathbf{f}(\mathbf{x}^j)$ , remember the solution that should be eliminated, and go to step 6.

5. If for all  $i = 1, 2, \dots, \ell$ ,
$$f_i(\mathbf{x}^j) \geq f_i(\mathbf{x}^{*r}),$$

then go to step 2.

6. Set  $r := r + 1$ . If  $r \leq R$ , go to step 4.



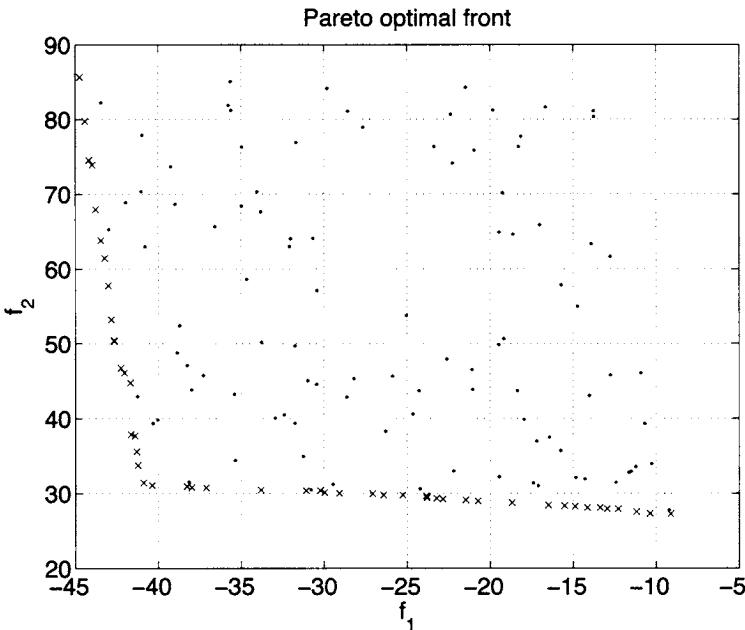
**Figure 24.5** Pareto optimal points in the decision space along with the level sets of the objective functions  $f_1$  and  $f_2$ .

7. If  $q \neq 0$ , remove from the candidate Pareto set the solutions that are eliminated in step 4, add solution  $\mathbf{x}^j$  as a new candidate Pareto solution, and go to step 2.
8. Set  $R := R + 1$ ,  $\mathbf{x}^{*R} := \mathbf{x}^j$ ,  $\mathbf{f}^{*R} := \mathbf{f}(\mathbf{x}^j)$ , and go to step 2.

**Example 24.2** We apply the algorithm above to generate the Pareto front for the multiobjective optimization problem

$$\begin{aligned} & \text{minimize} && \begin{bmatrix} -(x_1^2 + x_2) \\ x_1 + x_2^2 \end{bmatrix} \\ & \text{subject to} && 2 \leq x_1 \leq 6 \\ & && 5 \leq x_2 \leq 9. \end{aligned}$$

We performed 100 iterations. At each iteration we randomly generated 50 feasible points. Then we applied the algorithm above to extract from this set of feasible points candidate Pareto optimal solutions. In Figure 24.5 we show Pareto optimal points obtained after 100 iterations of the algorithm. We also show level sets of the objective functions in the  $(x_1, x_2)$ -space. In Figure 24.6 we show the Pareto front in the objective function space after 100



**Figure 24.6** Pareto front for the problem of Example 24.2. Also marked are the objective vector values for the remaining candidate points generated in the last iteration.

iterations of the algorithm. The Pareto optimal points are marked with x's. The points marked with '-'s are the candidate solutions generated randomly at the beginning of the last iteration of the algorithm. ■

We have described a simple approach to computing the Pareto front. Alternative methods include those that apply genetic algorithms to solving multiobjective optimization problems, as discussed in Deb [37], Coello Coello et al. [31], and Osyczka [98].

#### 24.4 From Multiobjective to Single-Objective Optimization

In some cases it is possible to deal with a multiobjective optimization problem by converting the problem into a single-objective optimization problem, so that standard optimization methods can be brought to bear. Here, we discuss four techniques to convert a multiobjective problem to a single-objective problem. We assume throughout that an objective function vector  $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), \dots, f_\ell(\mathbf{x})]^\top$  is given.

The first method is to form a single objective function by taking a linear combination, with positive coefficients, of the components of the objective

function vector. Equivalently, we form a convex combination of the components of the objective function vector. In other words, we use

$$f(\mathbf{x}) = \mathbf{c}^\top \mathbf{f}(\mathbf{x})$$

as the (single) objective function, where  $\mathbf{c}$  is a vector of positive components. This method is also called the *weighted-sum method*, where the coefficients of the linear combination (i.e., the components of  $\mathbf{c}$ ) are called *weights*. These weights reflect the relative importance of the individual components in the objective vector. Of course, it might be difficult to determine suitable weight values.

A second method is to form a single objective function by taking the maximum of the components of the objective vector:

$$f(\mathbf{x}) = \max\{f_1(\mathbf{x}), \dots, f_\ell(\mathbf{x})\}.$$

In other words, we convert the multiobjective minimization problem into one of minimizing the maximum of the components. For this reason, it is also called the *minimax method*. Note that this method applies to situations where the components of the objective vector are comparable or compatible, in the sense that they are in the same units (e.g., they are all lengths measured in meters, or masses in kilograms). A limitation of this method is that the resulting single objective function might not be differentiable, thereby precluding the use of optimization methods that rely on differentiability (e.g., gradient algorithms). However, as we show in the following, a minimax problem with linear objective vector components and linear constraints can be reduced to a linear programming problem.

**Example 24.3** Given vectors  $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^n$  and scalars  $u_1, \dots, u_p$ , consider the minimax problem

$$\begin{aligned} & \text{minimize} && \max\{\mathbf{v}_1^\top \mathbf{x} + u_1, \dots, \mathbf{v}_p^\top \mathbf{x} + u_p\} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{b} \in \mathbb{R}^m$ . Call this problem P1.

- a. Consider the optimization problem

$$\begin{aligned} & \text{minimize} && y \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && y \geq \mathbf{v}_i^\top \mathbf{x} + u_i, \quad i = 1, \dots, p, \end{aligned}$$

where the decision variable is the vector  $[\mathbf{x}^\top, y]^\top$ . Call this problem P2. Show that  $\mathbf{x}^*$  solves P1 if and only if  $[\mathbf{x}^*^\top, y^*]^\top$  with  $y^* = \max\{\mathbf{v}_1^\top \mathbf{x}^* + u_1, \dots, \mathbf{v}_p^\top \mathbf{x}^* + u_p\}$  solves P2.

*Hint:*  $y \geq \max\{a, b, c\}$  if and only if  $y \geq a$ ,  $y \geq b$ , and  $y \geq c$ .

- b. Use part a to derive a linear programming problem

$$\begin{aligned} & \text{minimize } \hat{\mathbf{c}}^\top \mathbf{z} \\ & \text{subject to } \hat{\mathbf{A}}\mathbf{z} \leq \hat{\mathbf{b}} \end{aligned}$$

that is equivalent to P1 (by “equivalent” we mean that the solution to one gives us the solution to the other). Explain how a solution to the linear programming problem above gives a solution to P1.

*Solution:*

- a. First suppose that  $\mathbf{x}^*$  is optimal in P1. Let  $y^* = \max\{\mathbf{v}_1^\top \mathbf{x}^* + u_1, \dots, \mathbf{v}_p^\top \mathbf{x}^* + u_p\}$ . Then,  $[\mathbf{x}^{*\top}, y^*]^\top$  is feasible in P2. Let  $[\mathbf{x}^\top, y]^\top$  be any feasible point in P2. Then (by the hint)

$$y \geq \max\{\mathbf{v}_1^\top \mathbf{x} + u_1, \dots, \mathbf{v}_p^\top \mathbf{x} + u_p\}.$$

Moreover,  $\mathbf{x}$  is feasible in P1, and hence

$$\begin{aligned} y &\geq \max\{\mathbf{v}_1^\top \mathbf{x} + u_1, \dots, \mathbf{v}_p^\top \mathbf{x} + u_p\} \\ &\geq \max\{\mathbf{v}_1^\top \mathbf{x}^* + u_1, \dots, \mathbf{v}_p^\top \mathbf{x}^* + u_p\} \\ &= y^*. \end{aligned}$$

Hence,  $[\mathbf{x}^{*\top}, y^*]^\top$  is optimal in the linear programming problem.

To prove the converse, suppose that  $\mathbf{x}^*$  is not optimal in P1. Then, there is some  $\mathbf{x}'$  that is feasible in P1 such that

$$\begin{aligned} y' &= \max\{\mathbf{v}_1^\top \mathbf{x}' + u_1, \dots, \mathbf{v}_p^\top \mathbf{x}' + u_p\} \\ &< \max\{\mathbf{v}_1^\top \mathbf{x}^* + u_1, \dots, \mathbf{v}_p^\top \mathbf{x}^* + u_p\} \\ &= y^*. \end{aligned}$$

But  $[\mathbf{x}'^\top, y']^\top$  is evidently feasible in P2, and has objective function value ( $y'$ ) that is lower than that of  $[\mathbf{x}^{*\top}, y^*]^\top$ . Hence,  $[\mathbf{x}^{*\top}, y^*]^\top$  is not optimal in P2.

- b. Define

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix}, \quad \hat{\mathbf{c}} = \begin{bmatrix} \mathbf{0} \\ 1 \end{bmatrix}, \quad \hat{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & 0 \\ \mathbf{v}_1^\top & -1 \\ \vdots & \vdots \\ \mathbf{v}_p^\top & -1 \end{bmatrix}, \quad \hat{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ -u_1 \\ \vdots \\ -u_p \end{bmatrix}.$$

Then the equivalent problem can be written as

$$\begin{aligned} & \text{minimize } \hat{\mathbf{c}}^\top \mathbf{z} \\ & \text{subject to } \hat{\mathbf{A}}\mathbf{z} \leq \hat{\mathbf{b}}. \end{aligned}$$

By part a, if we obtain a solution to this linear programming problem, then the first  $n$  components form a solution to the original minimax problem. ■

A third method to convert a multiobjective problem to a single-objective problem, assuming that the components of the objective vector are nonnegative, is to form a single objective function by taking the  $p$ -norm of the objective vector:

$$f(\mathbf{x}) = \|\mathbf{f}(\mathbf{x})\|_p.$$

The minimax method can be viewed as a special case of this method, with  $p = \infty$ . The weighted-sum method with uniform weights can be viewed as this method with  $p = 1$ . To make the objective function differentiable in the case where  $p$  is finite (so that we can apply gradient methods, for example), we replace it by its  $p$ th power:

$$f(\mathbf{x}) = \|\mathbf{f}(\mathbf{x})\|_p^p = (f_1(\mathbf{x}))^p + \cdots + (f_\ell(\mathbf{x}))^p.$$

A fourth method is to minimize one of the components of the objective vector subject to constraints on the other components. For example, given  $\mathbf{f}$ , we solve

$$\begin{aligned} & \text{minimize} && f_1(\mathbf{x}) \\ & \text{subject to} && f_2(\mathbf{x}) \leq b_2, \\ & && \vdots \\ & && f_\ell(\mathbf{x}) \leq b_\ell, \end{aligned}$$

where  $b_2, \dots, b_\ell$  are given constants that reflect satisfactory values for the objectives  $f_2, \dots, f_\ell$ , respectively. Of course, this approach is suitable only in situations where these satisfactory values can be determined.

## 24.5 Uncertain Linear Programming Problems

In this section we show how multiobjective optimization methods can be used to solve linear programming problems with uncertain coefficients, including uncertain constraints and uncertain objective functions.

### Uncertain Constraints

We first consider a generalization of linear programming to problems with uncertain constraints. Our exposition is based on a discussion of fuzzy linear programming by Wang [131, Chapter 30]. We consider the following general

linear programming problem:

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \quad \mathbf{x} \geq 0. \end{aligned}$$

We can represent the constraints in the form

$$(\mathbf{A}\mathbf{x})_i \leq b_i, \quad i = 1, 2, \dots, m.$$

Suppose that the constraints' bounds are uncertain in the sense that they can vary within given tolerance values and can be represented as

$$(\mathbf{A}\mathbf{x})_i \leq b_i + \theta t_i, \quad i = 1, 2, \dots, m,$$

where  $\theta \in [0, 1]$  and  $t_i > 0$ ,  $i = 1, 2, \dots, m$ .

We now discuss a method to solve the problem above. First, solve the following two linear programming problems:

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } (\mathbf{A}\mathbf{x})_i \leq b_i, \quad i = 1, 2, \dots, m \\ & \quad \mathbf{x} \geq 0 \end{aligned}$$

and

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } (\mathbf{A}\mathbf{x})_i \leq b_i + t_i, \quad i = 1, 2, \dots, m \\ & \quad \mathbf{x} \geq 0. \end{aligned}$$

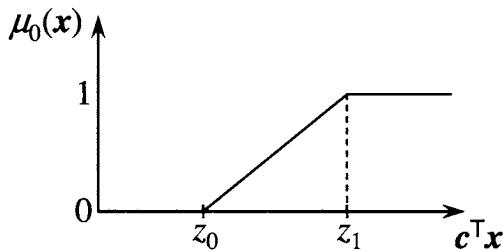
Denote the solution to the two programs as  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(0)}$ , respectively, and let  $z_1 = \mathbf{c}^\top \mathbf{x}^{(1)}$  and  $z_0 = \mathbf{c}^\top \mathbf{x}^{(0)}$ . Using these definitions, we construct a function that characterizes the “degree of the penalty” associated with the uncertain constraints in the linear programming problem

$$\mu_0(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{c}^\top \mathbf{x} < z_0 \\ \frac{\mathbf{c}^\top \mathbf{x} - z_0}{z_1 - z_0} & \text{if } z_0 \leq \mathbf{c}^\top \mathbf{x} \leq z_1 \\ 1 & \text{if } \mathbf{c}^\top \mathbf{x} > z_1. \end{cases}$$

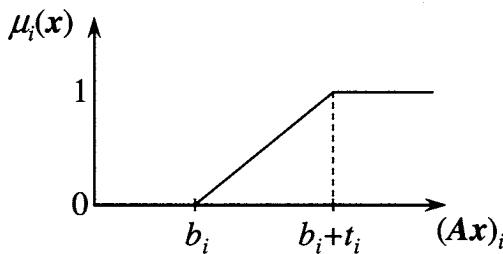
A plot of this function is given in Figure 24.7. Note that when  $\mathbf{c}^\top \mathbf{x} \leq z_0$ , then  $\mu_0(\mathbf{x}) = 0$ , which represents minimum degree of penalty. On the other hand, when  $\mathbf{c}^\top \mathbf{x} \geq z_1$ , then  $\mu_0(\mathbf{x}) = 1$ , and we have a maximum degree of penalty. When  $z_0 \leq \mathbf{c}^\top \mathbf{x} \leq z_1$ , the degree of penalty varies from 0 to 1.

Next, we introduce a function that describes the degree of penalty for violating the  $i$ th constraint:

$$\mu_i(\mathbf{x}) = \begin{cases} 0 & \text{if } (\mathbf{A}\mathbf{x})_i - b_i < 0 \\ \frac{(\mathbf{A}\mathbf{x})_i - b_i}{t_i} & \text{if } 0 \leq (\mathbf{A}\mathbf{x})_i - b_i \leq t_i \\ 1 & \text{if } (\mathbf{A}\mathbf{x})_i - b_i > t_i. \end{cases}$$



**Figure 24.7** Plot of the function  $\mu_0(\mathbf{x})$ .



**Figure 24.8** Plot of the function  $\mu_i(\mathbf{x})$ .

A plot of this function is shown in Figure 24.8.

Using the definitions above we can reformulate the original linear programming problem as a multiobjective optimization problem, with the goal of minimizing the functions that penalize constraint violations:

$$\begin{aligned} & \text{minimize} && \begin{bmatrix} \mu_0(\mathbf{x}) \\ \mu_1(\mathbf{x}) \\ \vdots \\ \mu_m(\mathbf{x}) \end{bmatrix} \\ & \text{subject to} && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

We can employ the minimax method to solve the multiobjective optimization problem as a single-objective problem

$$\begin{aligned} & \text{minimize} && \max \{\mu_0(\mathbf{x}), \mu_1(\mathbf{x}), \dots, \mu_m(\mathbf{x})\} \\ & \text{subject to} && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

As shown in Example 24.3, the problem above can be stated equivalently as

$$\begin{aligned} & \text{minimize } \theta \\ & \text{subject to } \mu_0(\mathbf{x}) \leq \theta \\ & \quad \mu_i(\mathbf{x}) \leq \theta, \quad i = 1, 2, \dots, m \\ & \quad \theta \in [0, 1], \quad \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Using now the definitions of  $\mu_0$  and  $\mu_i$ ,  $i = 1, \dots, m$ , we restate the optimization problem above as

$$\begin{aligned} & \text{minimize } \theta \\ & \text{subject to } \mathbf{c}^\top \mathbf{x} \leq z_0 + \theta(z_1 - z_0) \\ & \quad (\mathbf{A}\mathbf{x})_i \leq b_i + \theta t_i, \quad i = 1, 2, \dots, m \\ & \quad \theta \in [0, 1], \quad \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

**Example 24.4** Consider the following linear programming problem:

$$\begin{aligned} & \text{minimize } -\frac{1}{2}x_1 - x_2 \\ & \text{subject to } x_1 + x_2 \leq 5 \\ & \quad x_2 \leq 3 \\ & \quad x_1 \geq 0, \quad x_2 \geq 0, \end{aligned}$$

where the tolerances are  $t_1 = 2$  and  $t_2 = 1$ .

- Solve the two linear programming problems to obtain  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(0)}$  using the data above. Then find  $z_1$  and  $z_0$ .
- Construct the equivalent optimization problem (involving  $\theta$ ) using the data above.
- Express the optimization problem as a linear programming problem in standard form.

*Solution:*

- We can solve these problems graphically to obtain

$$\mathbf{x}^{(1)} = [2, 3]^\top \quad \text{and} \quad \mathbf{x}^{(0)} = [3, 4]^\top.$$

Hence,

$$z_1 = \mathbf{c}^\top \mathbf{x}^{(1)} = -4 \quad \text{and} \quad z_0 = \mathbf{c}^\top \mathbf{x}^{(0)} = -5\frac{1}{2}.$$

b. The optimization problem has the form

$$\begin{aligned} & \text{minimize } \theta \\ & \text{subject to } \mu_0(\mathbf{x}) \leq \theta \\ & \quad \mu_1(\mathbf{x}) \leq \theta \\ & \quad \mu_2(\mathbf{x}) \leq \theta \\ & \quad \theta \in [0, 1], \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \mu_0(\mathbf{x}) &= \begin{cases} 0 & \text{if } -\frac{1}{2}x_1 - x_2 < -5\frac{1}{2} \\ \frac{-\frac{1}{2}x_1 - x_2 + 5\frac{1}{2}}{3/2} & \text{if } -5\frac{1}{2} \leq -\frac{1}{2}x_1 - x_2 \leq -4 \\ 1 & \text{if } -\frac{1}{2}x_1 - x_2 > -4, \end{cases} \\ \mu_1(\mathbf{x}) &= \begin{cases} 0 & \text{if } x_1 + x_2 - 5 < 0 \\ \frac{x_1 + x_2 - 5}{2} & \text{if } 0 \leq x_1 + x_2 - 5 \leq 2 \\ 1 & \text{if } x_1 + x_2 - 5 > 2, \end{cases} \\ \mu_2(\mathbf{x}) &= \begin{cases} 0 & \text{if } x_2 - 3 < 0 \\ x_2 - 3 & \text{if } 0 \leq x_2 - 3 \leq 1 \\ 1 & \text{if } x_2 - 3 > 1. \end{cases} \end{aligned}$$

c. We have

$$\begin{aligned} & \text{minimize } \theta \\ & \text{subject to } \mathbf{c}^\top \mathbf{x} \leq z_0 + \theta(z_1 - z_0) \\ & \quad (\mathbf{Ax})_i \leq b_i + (1 - \theta)t_i, \quad i = 1, 2 \\ & \quad \theta \in [0, 1], \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Using our data, we obtain

$$\begin{aligned} & \text{minimize } \theta \\ & \text{subject to } \frac{1}{2}x_1 + x_2 \geq 5\frac{1}{2} - \frac{3}{2}\theta \\ & \quad x_1 + x_2 \leq 5 + 2\theta \\ & \quad x_2 \leq 3 + \theta \\ & \quad \theta \in [0, 1], \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

Write  $x_3 = \theta$ . Then, the above problem can be represented as

$$\begin{aligned} & \text{minimize } x_3 \\ & \text{subject to } x_1 + 2x_2 + 3x_3 \geq 11 \\ & \quad x_1 + x_2 - 2x_3 \leq 5 \\ & \quad x_2 - x_3 \leq 3 \\ & \quad x_3 \leq 1 \\ & \quad x_i \geq 0, \quad i = 1, 2, 3. \end{aligned}$$

The above linear program expressed in the form of a linear programming problem in standard form is

$$\begin{aligned} & \text{minimize } x_3 \\ & \text{subject to } x_1 + 2x_2 + 3x_3 - x_4 = 11 \\ & \quad x_1 + x_2 - 2x_3 + x_5 = 5 \\ & \quad x_2 - x_3 + x_6 = 3 \\ & \quad x_3 + x_7 = 1 \\ & \quad x_i \geq 0, \quad i = 1, 2, \dots, 7. \end{aligned}$$

■

### Uncertain Objective Function Coefficients

We now consider a linear programming problem with uncertain objective function coefficients. We assume that uncertainties of the objective coefficients are modeled by the following triangular function:

$$\mu(x; a, b, c) = \begin{cases} 0 & \text{if } x < a \\ (x-a)/(b-a) & \text{if } a \leq x < b \\ (c-x)/(c-b) & \text{if } b \leq x \leq c \\ 0 & \text{if } x > c. \end{cases}$$

A plot of this function for  $a = 1$ ,  $b = 2$ , and  $c = 6$  is shown in Figure 24.9. In other words, the uncertain objective coefficients will be represented by the triangular functions of the form given above. Following Wang [131, p. 386], we use the notation  $\tilde{c}_i = (c_i^-, c_i^0, c_i^+)$  to denote the uncertain coefficient  $c_i$  represented by the triangular function  $\mu(x; c_i^-, c_i^0, c_i^+)$ . Then the linear programming problem

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{x} \\ & \text{subject to } \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \quad \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

becomes

$$\begin{aligned} & \text{minimize} && \begin{bmatrix} \mathbf{c}^-\mathbf{x} \\ \mathbf{c}^0\mathbf{x} \\ \mathbf{c}^+\mathbf{x} \end{bmatrix} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where

$$\mathbf{c}^- = [c_1^- \quad \dots \quad c_n^-], \quad \mathbf{c}^0 = [c_1^0 \quad \dots \quad c_n^0], \quad \mathbf{c}^+ = [c_1^+ \quad \dots \quad c_n^+].$$

This is a multiobjective optimization problem. Wang [131] suggests that instead of minimizing the three values  $\mathbf{c}^-\mathbf{x}$ ,  $\mathbf{c}^0\mathbf{x}$ , and  $\mathbf{c}^+\mathbf{x}$  simultaneously, the center,  $\mathbf{c}^0\mathbf{x}$ , be minimized; the left leg,  $(\mathbf{c}^0 - \mathbf{c}^-)\mathbf{x}$ , be maximized; and the right leg,  $(\mathbf{c}^+ - \mathbf{c}^0)\mathbf{x}$ , be minimized. This results in pushing the triangular functions to the left in the minimization process. Thus, the multiobjective optimization problem above can be changed to the following multiobjective optimization problem:

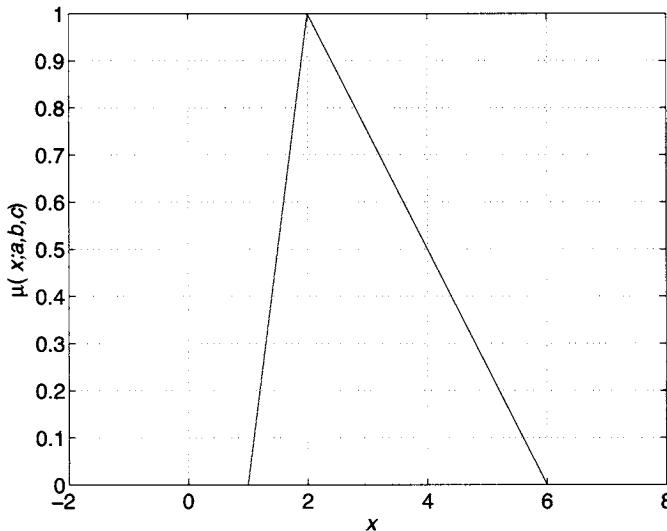
$$\begin{aligned} & \text{minimize} && \begin{bmatrix} -(\mathbf{c}^0 - \mathbf{c}^-)\mathbf{x} \\ \mathbf{c}^0\mathbf{x} \\ (\mathbf{c}^+ - \mathbf{c}^0)\mathbf{x} \end{bmatrix} \\ & \text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

### Uncertain Constraint Coefficients

We may be faced with solving a linear programming problem with uncertain constraint coefficients. In this case the coefficients of the constraint matrix  $\mathbf{A}$  would be represented by triangular functions of the form given in the preceding section. That is, the coefficient  $a_{ij}$  of the constraint matrix  $\mathbf{A}$  would be modeled by the function  $\tilde{a}_{ij} = \mu(x; a_{ij}^-, a_{ij}^0, a_{ij}^+)$ . Then, the linear programming problem with uncertain constraint coefficients would take the form

$$\begin{aligned} & \text{minimize} && \mathbf{c}^\top \mathbf{x} \\ & \text{subject to} && \begin{bmatrix} \mathbf{A}^-\mathbf{x} \\ \mathbf{A}^0\mathbf{x} \\ \mathbf{A}^+\mathbf{x} \end{bmatrix} \leq \begin{bmatrix} \mathbf{b} \\ \mathbf{b} \\ \mathbf{b} \end{bmatrix} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{A}^- = [a_{ij}^-]$ ,  $\mathbf{A}^0 = [a_{ij}^0]$ , and  $\mathbf{A}^+ = [a_{ij}^+]$ .



**Figure 24.9** Plot of the triangular function  $\mu(x; a, b, c)$  for  $a = 1$ ,  $b = 2$ , and  $c = 6$ .

### General Uncertainties

Finally, we may be faced with solving an uncertain linear programming problem that is a combination of the basic uncertain linear programming problems discussed above. For example, suppose that we are asked to solve the following quite general uncertain linear programming problem:

$$\begin{aligned} & \text{minimize} && \tilde{\mathbf{c}}^\top \mathbf{x} \\ & \text{subject to} && \tilde{\mathbf{A}}\mathbf{x} \leq \tilde{\mathbf{b}} \\ & && \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where the tilde symbols refer to the uncertain data; that is, we have

$$\tilde{\mathbf{c}} = (\mathbf{c}^-, \mathbf{c}^0, \mathbf{c}^+), \quad \tilde{\mathbf{A}} = (\mathbf{A}^-, \mathbf{A}^0, \mathbf{A}^+), \quad \tilde{\mathbf{b}} = (\mathbf{b}^-, \mathbf{b}^0, \mathbf{b}^+).$$

We can represent the uncertain linear programming problem above as a multiobjective optimization problem of the form

$$\begin{aligned} & \text{minimize} && \begin{bmatrix} -(\mathbf{c}^0 - \mathbf{c}^-)\mathbf{x} \\ \mathbf{c}^0\mathbf{x} \\ (\mathbf{c}^+ - \mathbf{c}^0)\mathbf{x} \end{bmatrix} \\ & \text{subject to} && \begin{bmatrix} \mathbf{A}^-\mathbf{x} \\ \mathbf{A}^0\mathbf{x} \\ \mathbf{A}^+\mathbf{x} \end{bmatrix} \leq \begin{bmatrix} \mathbf{b}^- \\ \mathbf{b}^0 \\ \mathbf{b}^+ \end{bmatrix} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

## EXERCISES

**24.1** Write a MATLAB program that implements the algorithm for generating a Pareto front, and test it on the problem in Example 24.1.

**24.2** Consider the multiobjective problem

$$\begin{aligned} & \text{minimize } \mathbf{f}(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$ .

- a. Suppose that we solve the single-objective problem

$$\begin{aligned} & \text{minimize } \mathbf{c}^\top \mathbf{f}(\mathbf{x}) \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $\mathbf{c} \in \mathbb{R}^n$ ,  $\mathbf{c} > \mathbf{0}$  (i.e., we use the weighted-sum approach). Show that if  $\mathbf{x}^*$  is a global minimizer for the single-objective problem above, then  $\mathbf{x}^*$  is a Pareto minimizer for the given multiobjective problem. Then show that it is not necessarily the case that if  $\mathbf{x}^*$  is a Pareto minimizer for the multiobjective problem, then there exists a  $\mathbf{c} > \mathbf{0}$  such that  $\mathbf{x}^*$  is a global minimizer for the single-objective (weighted-sum) problem.

- b. Assuming that for all  $\mathbf{x} \in \Omega$ ,  $\mathbf{f}(\mathbf{x}) \geq \mathbf{0}$ , suppose that we solve the single-objective problem

$$\begin{aligned} & \text{minimize } (f_1(\mathbf{x}))^p + \cdots + (f_\ell(\mathbf{x}))^p \\ & \text{subject to } \mathbf{x} \in \Omega, \end{aligned}$$

where  $p \in \mathbb{R}$ ,  $p > 0$  (i.e., we use the minimum-norm approach). Show that if  $\mathbf{x}^*$  is a global minimizer for the single-objective problem above, then  $\mathbf{x}^*$  is a Pareto minimizer for the given multiobjective problem. Then show that it is not necessarily the case that if  $\mathbf{x}^*$  is a Pareto minimizer for the multiobjective problem, then there exists a  $p > 0$  such that  $\mathbf{x}^*$  is a global minimizer for the single-objective (minimum-norm) problem.

- c. Suppose that we solve the single-objective problem

$$\begin{aligned} & \text{minimize } \max\{f_1(\mathbf{x}), \dots, f_\ell(\mathbf{x})\} \\ & \text{subject to } \mathbf{x} \in \Omega \end{aligned}$$

(i.e., we use the minimax approach). Show that it is not necessarily the case that if  $\mathbf{x}^*$  is a Pareto minimizer for the given multiobjective problem, then  $\mathbf{x}^*$  is a global minimizer for the single-objective (minimax) problem. Then show that it also is not necessarily the case that if  $\mathbf{x}^*$  is a global

minimizer for the single-objective problem, then  $\mathbf{x}^*$  is a Pareto minimizer for the multiobjective problem.

**24.3** Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$  be given. Consider the following multiobjective problem with equality constraints:

$$\begin{aligned} & \text{minimize} && \mathbf{f}(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \Omega. \end{aligned}$$

Suppose that  $\mathbf{f} \in \mathcal{C}^1$ , all the components of  $\mathbf{f}$  are convex, and  $\Omega$  is convex. Suppose that there exists  $\mathbf{x}^*$  and  $\mathbf{c}^* > \mathbf{0}$  such that for any feasible direction  $\mathbf{d}$  at  $\mathbf{x}^*$ , we have

$$\mathbf{c}^{*\top} D\mathbf{f}(\mathbf{x}^*) \mathbf{d} \geq 0.$$

Show that  $\mathbf{x}^*$  is a Pareto minimizer.

**24.4** Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$  and  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be given. Consider the following multiobjective problem with equality constraints:

$$\begin{aligned} & \text{minimize} && \mathbf{f}(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0}. \end{aligned}$$

Suppose that  $\mathbf{f}, \mathbf{h} \in \mathcal{C}^1$ , all the components of  $\mathbf{f}$  are convex, and the constraint set is convex. Show that if there exists  $\mathbf{x}^*$ ,  $\mathbf{c}^* > \mathbf{0}$ , and  $\boldsymbol{\lambda}^*$  such that

$$\begin{aligned} \mathbf{c}^{*\top} D\mathbf{f}(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} Dh(\mathbf{x}^*) &= \mathbf{0}^\top \\ \mathbf{h}(\mathbf{x}^*) &= \mathbf{0}, \end{aligned}$$

then  $\mathbf{x}^*$  is a Pareto minimizer. We can think of the above as a Lagrange condition for the constrained multiobjective function.

**24.5** Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$  and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$  be given. Consider the following multiobjective problem with inequality constraints:

$$\begin{aligned} & \text{minimize} && \mathbf{f}(\mathbf{x}) \\ & \text{subject to} && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

Suppose that  $\mathbf{f}, \mathbf{g} \in \mathcal{C}^1$ , all the components of  $\mathbf{f}$  are convex, and the constraint set is convex. Show that if there exists  $\mathbf{x}^*$ ,  $\mathbf{c}^* > \mathbf{0}$ , and  $\boldsymbol{\mu}^*$  such that

$$\begin{aligned} \boldsymbol{\mu}^* &\geq \mathbf{0}, \\ \mathbf{c}^{*\top} D\mathbf{f}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} Dg(\mathbf{x}^*) &= \mathbf{0}^\top, \\ \boldsymbol{\mu}^{*\top} g(\mathbf{x}^*) &= 0, \\ g(\mathbf{x}^*) &\leq \mathbf{0}, \end{aligned}$$

then  $\mathbf{x}^*$  is a Pareto minimizer. We can think of the above as a KKT condition for the constrained multiobjective function.

**24.6** Let  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^\ell$ ,  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , and  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$  be given. Consider the general constrained multiobjective problem

$$\begin{aligned} & \text{minimize} && \mathbf{f}(\mathbf{x}) \\ & \text{subject to} && \mathbf{h}(\mathbf{x}) = \mathbf{0} \\ & && \mathbf{g}(\mathbf{x}) \leq \mathbf{0}. \end{aligned}$$

Suppose that  $\mathbf{f}, \mathbf{g}, \mathbf{h} \in \mathcal{C}^1$ , all the components of  $\mathbf{f}$  are convex, and the constraint set is convex. Show that if there exists  $\mathbf{x}^*$ ,  $\mathbf{c}^* > \mathbf{0}$ ,  $\boldsymbol{\lambda}^*$ , and  $\boldsymbol{\mu}^*$  such that

$$\begin{aligned} & \boldsymbol{\mu}^* \geq \mathbf{0}, \\ & \mathbf{c}^{*\top} D\mathbf{f}(\mathbf{x}^*) + \boldsymbol{\lambda}^{*\top} D\mathbf{h}(\mathbf{x}^*) + \boldsymbol{\mu}^{*\top} D\mathbf{g}(\mathbf{x}^*) = \mathbf{0}^\top, \\ & \boldsymbol{\mu}^{*\top} \mathbf{g}(\mathbf{x}^*) = 0, \\ & \mathbf{h}(\mathbf{x}^*) = \mathbf{0}, \\ & \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}, \end{aligned}$$

then  $\mathbf{x}^*$  is a Pareto minimizer.

**24.7** Let  $f_1 : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f_1, f_2 \in \mathcal{C}^1$ . Consider the minimax problem

$$\text{minimize } \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\}.$$

Show that if  $\mathbf{x}^*$  is a local minimizer, then there exist  $\mu_1^*, \mu_2^* \in \mathbb{R}$  such that

$$\mu_1^*, \mu_2^* \geq 0, \quad \mu_1^* \nabla f_1(\mathbf{x}^*) + \mu_2^* \nabla f_2(\mathbf{x}^*) = \mathbf{0}, \quad \mu_1^* + \mu_2^* = 1,$$

and  $\mu_i^* = 0$  if  $f_i(\mathbf{x}^*) < \max\{f_1(\mathbf{x}^*), f_2(\mathbf{x}^*)\}$ .

*Hint:* Consider the following problem: minimize  $z$  subject to  $z \geq f_i(\mathbf{x})$ ,  $i = 1, 2$ .

## REFERENCES

---

1. J. S. Arora, *Introduction to Optimum Design*. New York: McGraw-Hill Book Co., 1989.
2. R. G. Bartle, *The Elements of Real Analysis*, 2nd ed. New York: Wiley, 1976.
3. M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*, 2nd ed. New York: Wiley, 1993.
4. A. Bhaya and E. Kaszkurewicz, *Control Perspectives on Numerical Algorithms and Matrix Problems*. Philadelphia: Society for Industrial and Applied Mathematics, 2006.
5. B. Beliczynski, A. Dzielinski, M. Iwanowski, and B. Ribeiro, Eds., *Adaptive and Natural Computing Algorithms*, vol. 4431 of *Lecture Notes in Computer Science*. Berlin: Springer, 2007.
6. A. Ben-Israel and T. N. E. Greville, *Generalized Inverses: Theory and Applications*. New York: Wiley-Interscience, 1974.
7. L. D. Berkovitz, *Convexity and Optimization in  $\mathbb{R}^n$* . Hoboken, NJ: Wiley, 2002.
8. C. C. Berresford, A. M. Rockett, and J. C. Stevenson, “Khachiyan’s algorithm, Part 1: A new solution to linear programming problems,” *Byte*, vol. 5, no. 8, pp. 198–208, Aug. 1980.
9. C. C. Berresford, A. M. Rockett, and J. C. Stevenson, “Khachiyan’s algorithm, Part 2: Problems with the algorithm,” *Byte*, vol. 5, no. 9, pp. 242–255, Sept. 1980.

10. D. P. Bertsekas, "Necessary and sufficient conditions for a penalty method to be exact," *Mathematical Programming*, vol. 9, no. 1, pp. 87–99, Aug. 1975.
11. D. P. Bertsekas, *Nonlinear Programming: 2nd ed.* Belmont, MA: Athena Scientific, 1999.
12. D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Belmont, MA: Athena Scientific, 1997.
13. K. G. Binmore, *Calculus*. Cambridge, England: Cambridge University Press, 1986.
14. R. G. Bland, D. Goldfarb, and M. J. Todd, "The ellipsoid method: A survey," *Operations Research*, vol. 29, pp. 1039–1091, 1981.
15. V. G. Boltyanskii, *Mathematical Methods of Optimal Control*. New York: Holt, Rinehart and Winston, 1971.
16. S. Boyd, L. El Ghaoui, E. Feron, and V. Balakrishnan, *Linear Matrix Inequalities in System and Control Theory*. Philadelphia, PA: SIAM, 1994.
17. R. P. Brent, *Algorithms for Minimization without Derivatives*. Englewood Cliffs, NJ: Prentice Hall, 1973.
18. L. Brickman, *Mathematical Introduction to Linear Programming and Game Theory*. New York: Springer-Verlag, 1989.
19. C. G. Broyden, "Quasi-Newton methods," in *Optimization Methods in Electronics and Communications* (K. W. Cattermole and J. J. O'Reilly, Eds.), vol. 1 of *Mathematical Topics in Telecommunications*. New York: Wiley, 1984, pp. 105–110,
20. A. E. Bryson and Y.-C. Ho, *Applied Optimal Control: Optimization, Estimation, and Control*, rev. print. Washington, DC: Hemisphere Publishing Corporation, 1975.
21. B. D. Bunday, *Basic Optimization Methods*. London: Edward Arnold, 1984.
22. J. Campbell, *The Improbable Machine*. New York: Simon and Schuster, 1989.
23. S. L. Campbell and C. D. Meyer, Jr., *Generalized Inverses of Linear Transformations*. New York: Dover Publications, 1991.
24. E. K. P. Chong and B. E. Brewington, "Distributed communications resource management for tracking and surveillance networks," in *Proceedings of the Conference on Signal and Data Processing of Small Targets 2005* (SPIE Vol. 5913), part of the *SPIE Symposium on Optics & Photonics*, San Diego, California, July 31–Aug. 4, 2005, pp. 280–291.
25. E. K. P. Chong and B. E. Brewington, "Decentralized rate control for tracking and surveillance networks," *Ad Hoc Networks*, special issue on *Recent Advances in Wireless Sensor Networks*, vol. 5, no. 6, pp. 910–928, Aug. 2007.
26. E. K. P. Chong, S. Hui, and S. H. Žak, "An analysis of a class of neural networks for solving linear programming problems," *IEEE Transactions on Automatic Control*, special section on *Neural Networks in Control, Identification, and Decision Making*, vol. 44, no. 11, pp. 1995–2006, Nov. 1999.
27. E. K. P. Chong and S. H. Žak, "Single-dimensional search methods," in *Wiley Encyclopedia of Operations Research and Management Science*, 2011, ISBN: 978-0-470-40063-0.

28. A. Cichocki and R. Unbehauen, *Neural Networks for Optimization and Signal Processing*. Chichester, England: Wiley, 1993.
29. M. Clerc, "The swarm and the queen: Towards a deterministic and adaptive particle swarm optimization," in *Proceedings of the Congress of Evolutionary Computation*, Washington, DC, July 1999, pp. 1951–1957.
30. M. Clerc and J. Kennedy, "The particle swarm: Explosion, stability and convergence in a multidimensional complex space," *IEEE Transactions on Evolutionary Computation*, vol. 6, pp. 58–73, Feb. 2002.
31. C. A. Coello Coello, D. A. Van Veldhuizen, and G. B. Lamont, *Evolutionary Algorithms for Solving Multi-Objective Problems*. New York: Kluwer Academic/Plenum Publishers, 2002.
32. S. D. Conte and C. de Boor, *Elementary Numerical Analysis: An Algorithmic Approach*, 3rd ed. New York: McGraw-Hill Book Co., 1980.
33. M. A. Dahleh and I. J. Diaz-Bobillo, *Control of Uncertain Systems: A Linear Programming Approach*. Upper Saddle River, NJ: Prentice Hall, 1995.
34. G. B. Dantzig, *Linear Programming and Extensions*. Princeton, NJ: Princeton University Press, 1963.
35. G. B. Dantzig and M. N. Thapa, *Linear Programming*, vol. 1, *Introduction*. New York: Springer-Verlag, 1997.
36. L. Davis, Ed., *Genetic Algorithms and Simulated Annealing*, Research Notes in Artificial Intelligence. London: Pitman, 1987.
37. K. Deb, *Multi-objective Optimization Using Evolutionary Algorithms*. Chichester, England: Wiley, 2001.
38. V. F. Dem'yanov and L. V. Vasil'ev, *Nondifferentiable Optimization*. New York: Optimization Software, Inc., Publications Division, 1985.
39. J. E. Dennis, Jr. and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Englewood Cliffs, NJ: Prentice Hall, 1983.
40. J. Dongarra and F. Sullivan, "The top 10 algorithms," *Computing in Science and Engineering*, pp. 22–23, Jan./Feb. 2000.
41. V. N. Faddeeva, *Computational Methods of Linear Algebra*. New York: Dover Publications, 1959.
42. S.-C. Fang and S. Puthenpura, *Linear Optimization and Extensions: Theory and Algorithms*. Englewood Cliffs, NJ: Prentice Hall, 1993.
43. R. Fletcher, *Practical Methods of Optimization*, 2nd ed. Chichester, England: Wiley, 1987.
44. F. R. Gantmacher, *The Theory of Matrices*, vol. 1. New York: Chelsea Publishing Co., 1959.
45. F. R. Gantmacher, *The Theory of Matrices*, 2nd ed. Moscow: Nauka, revised 1966. In Russian.
46. S. I. Gass, *An Illustrated Guide to Linear Programming*. New York: McGraw-Hill Book Co., 1970.

47. I. M. Gel'fand, *Lectures on Linear Algebra*. New York: Interscience Publishers, 1961.
48. S. Geman and D. Geman, "Stochastic relaxation, Gibbs distribution, and the Bayesian restoration of images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 6, pp. 721–741, 1984.
49. P. E. Gill and W. Murray, "Safeguarded steplength algorithms for optimization using descent methods," Tech. Rep. NPL NAC 37, National Physical Laboratory, Division of Numerical Analysis and Computing, Teddington, England, Aug. 1974.
50. P. E. Gill, W. Murray, M. A. Saunders, and M. H. Wright, "Two step-length algorithms for numerical optimization," Tech. Rep. SOL 79-25, Systems Optimization Laboratory, Department of Operations Research, Stanford University, Stanford, CA, Dec. 1979.
51. P. E. Gill, W. Murray, and M. H. Wright, *Practical Optimization*. London: Academic Press, 1981.
52. P. E. Gill, W. Murray, and M. H. Wright, *Numerical Linear Algebra and Optimization*. Redwood City, CA: Addison-Wesley, 1991.
53. G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed.. Baltimore, MD: The Johns Hopkins University Press, 1983.
54. R. E. Gomory, "Outline of an algorithm for integer solutions to linear programs," *Bulletin of the American Mathematical Society*, vol. 64, no. 5, pp. 275–278, Sep. 1958.
55. C. C. Gonzaga, "Path-following methods for linear programming," *SIAM Review*, vol. 34, no. 2, pp. 167–224, June 1992.
56. B. Hajek, "Cooling schedules for optimal annealing," *Mathematics of Operations Research*, vol. 13, no. 2, pp. 311–329, 1988.
57. J. Hannig, E. K. P. Chong, and S. R. Kulkarni, "Relative frequencies of generalized simulated annealing," *Mathematics of Operations Research*, vol. 31, no. 1, pp. 199–216, Feb. 2006.
58. R. L. Harvey, *Neural Network Principles*. Englewood Cliffs, NJ: Prentice Hall, 1994.
59. S. Haykin, *Neural Networks: A Comprehensive Foundation*, 2nd ed. Upper Saddle River, NJ: Prentice Hall, 1999.
60. J. Hertz, A. Krogh, and R. G. Palmer, *Introduction to the Theory of Neural Computation*, vol. 1 of *Santa Fe Institute Studies in the Sciences of Complexity*. Redwood City, CA: Addison-Wesley, 1991.
61. J. H. Holland, *Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence*. Cambridge, MA: MIT Press, 1992.
62. R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, England: Cambridge University Press, 1985.
63. A. S. Householder, *The Theory of Matrices in Numerical Analysis*. New York: Dover Publications, 1975.

64. S. Hui and S. H. Źak, "The Widrow-Hoff algorithm for McCulloch-Pitts type neurons," *IEEE Transactions on Neural Networks*, vol. 5, no. 6, pp. 924–929, Nov. 1994.
65. D. R. Hush and B. G. Horne, "Progress in supervised neural networks: What's new since Lippmann," *IEEE Signal Processing Magazine*, pp. 8–39, Jan. 1993.
66. S. Isaak and M. N. Manougian, *Basic Concepts of Linear Algebra*. New York: W. W. Norton & Co., 1976.
67. J.-S. R. Jang, C.-T. Sun, and E. Mizutani, *Neuro-Fuzzy and Soft Computing: A Computational Approach to Learning and Machine Intelligence*. Upper Saddle River, NJ: Prentice Hall, 1997.
68. W. E. Jenner, *Rudiments of Algebraic Geometry*. New York: Oxford University Press, 1963.
69. E. M. Johansson, F. U. Dowla, and D. M. Goodman, "Backpropagation learning for multi-layer feed-forward neural networks using the conjugate gradient method," *International Journal of Neural Systems*, vol. 2, no. 4, pp. 291–301, 1992.
70. S. Kaczmarz, "Approximate solution of systems of linear equations," *International Journal of Control*, vol. 57, no. 6, pp. 1269–1271, 1993. A reprint of the original paper: S. Kaczmarz, "Angenäherte Auflösung von Systemen linearer Gleichungen," *Bulletin International de l'Academie Polonaise des Sciences et des Lettres, Serie A*, pp. 355–357, 1937.
71. N. Karmarkar, "A new polynomial-time algorithm for linear programming," *Combinatorica*, vol. 4, no. 4, pp. 373–395, 1984.
72. M. F. Kelly, P. A. Parker, and R. N. Scott, "The application of neural networks to myoelectric signal analysis: A preliminary study," *IEEE Transactions on Biomedical Engineering*, vol. 37, no. 3, pp. 221–230, Mar. 1990.
73. J. Kennedy and R. C. Eberhart, with Y. Shi, *Swarm Intelligence*. San Francisco: Morgan Kaufmann, 2001.
74. L. G. Khachiyan, "A polynomial algorithm in linear programming," *Soviet Mathematics Doklady*, vol. 20, no. 1, pp. 191–194, 1979.
75. S. Kirkpatrick, C. D. Gelatt, Jr., and M. P. Vecchi, "Optimization by simulated annealing," *Science*, vol. 220, no. 4598, pp. 671–680, 1983.
76. V. Klee and G. J. Minty, "How good is the simplex algorithm?" in *Inequalities-III* (O. Shisha, Ed.), New York: Academic Press, 1972, pp. 159–175.
77. D. E. Knuth, *The Art of Computer Programming*, vol. 1, *Fundamental Algorithms*, 2nd ed. Reading, MA: Addison-Wesley, 1973.
78. L. Kolev, "Iterative algorithm for the minimum fuel and minimum amplitude problems for linear discrete systems," *International Journal of Control*, vol. 21, no. 5, pp. 779–784, 1975.
79. J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. Cambridge, MA: MIT Press, 1992.
80. T. Kozek, T. Roska, and L. O. Chua, "Genetic algorithm for CNN template learning," *IEEE Transactions on Circuits and Systems, I: Fundamental Theory and Applications*, vol. 40, no. 6, pp. 392–402, June 1993.

81. K. Kuratowski, *Introduction to Calculus*, 2nd ed., vol. 17 of *International Series of Monographs in Pure and Applied Mathematics*. Warsaw, Poland: Pergamon Press, 1969.
82. J. C. Lagarias, J. A. Reeds, M. H. Wright, and P. E. Wright, “Convergence properties of the Nelder-Mead simplex method in low dimensions,” *SIAM Journal on Optimization*, vol. 9, no. 1, pp. 112–147, 1998.
83. S. Lang, *Calculus of Several Variables*, 3rd ed. New York: Springer-Verlag, 1987.
84. J. M. Layton, *Multivariable Control Theory*. Stevenage, England: Peter Peregrinus on behalf of the Institution of Electrical Engineers, 1976.
85. E. B. Lee and L. Markus, *Foundations of Optimal Control Theory*. Malabar, FL: Robert E. Krieger Publishing Company, 1986.
86. G. Leitmann, *The Calculus of Variations and Optimal Control: An Introduction*. New York: Plenum Press, 1981.
87. D. G. Luenberger, *Optimization by Vector Space Methods*. New York: Wiley, 1969.
88. D. G. Luenberger and Y. Ye, *Linear and Nonlinear Programming*, 3rd ed. New York, NY: Springer Science + Business Media, 2008.
89. I. J. Maddox, *Elements of Functional Analysis*, 2nd ed. Cambridge, England: Cambridge University Press, 1988.
90. O. L. Mangasarian, *Nonlinear Programming*. New York: McGraw-Hill Book Co., 1969.
91. N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, H. Teller, and E. Teller, “Equation of state calculations by fast computing machines,” *Journal of Chemical Physics*, vol. 21, no. 6, pp. 1087–1092, 1953.
92. K. M. Miettinen, *Nonlinear Multiobjective Optimization*. Norwell, MA: Kluwer Academic Publishers, 1998.
93. S. A. Miller and E. K. P. Chong, “Flow-rate control for managing communications in tracking and surveillance networks,” in *Proceedings of the Conference on Signal and Data Processing of Small Targets 2007* (SPIE Vol. 6699), part of the *SPIE Symposium on Optics & Photonics*, San Diego, California, Aug. 26–30, 2007.
94. M. Mitchell, *An Introduction to Genetic Algorithms*. Cambridge, MA: MIT Press, 1996.
95. A. Mostowski and M. Stark, *Elements of Higher Algebra*. Warsaw, Poland: PWN—Polish Scientific Publishers, 1958.
96. S. G. Nash and A. Sofer, *Linear and Nonlinear Programming*. New York: McGraw-Hill Book Co., 1996.
97. J. A. Nelder and R. Mead, “A simplex method for function minimization,” *Computer Journal*, vol. 7, no. 4, pp. 308–313, 1965.
98. A. Osyczka, *Evolutionary Algorithms for Single and Multicriteria Design Optimization*. Heidelberg, Germany: Physica-Verlag, 2002.

99. D. H. Owens, *Multivariable and Optimal Systems*. London: Academic Press, 1981.
100. T. M. Ozan, *Applied Mathematical Programming for Production and Engineering Management*. Englewood Cliffs, NJ: Prentice Hall, 1986.
101. C. H. Papadimitriou and K. Steiglitz, *Combinatorial Optimization: Algorithms and Complexity*. Englewood Cliffs, NJ: Prentice Hall, 1982.
102. P. C. Parks, "S. Kaczmarz (1895–1939)," *International Journal of Control*, vol. 57, no. 6, pp. 1263–1267, 1993.
103. R. J. Patton and G. P. Liu, "Robust control design via eigenstructure assignment, genetic algorithms and gradient-based optimisation," *IEE Proceedings on Control Theory and Applications*, vol. 141, no. 3, pp. 202–208, May 1994.
104. A. L. Peressini, F. E. Sullivan, and J. J. Uhl, Jr., *The Mathematics of Nonlinear Programming*. New York: Springer-Verlag, 1988.
105. A. Pezeshki, L. L. Scharf, M. Lundberg, and E. K. P. Chong, "Constrained quadratic minimizations for signal processing and communications," in *Proceedings of the Joint 44th IEEE Conference on Decision and Control and European Control Conference (CDC-ECC'05)*, Seville, Spain, Dec. 12–15, 2005, pp. 7949–7953.
106. A. Pezeshki, L. L. Scharf, and E. K. P. Chong, "The geometry of linearly and quadratically constrained optimization problems for signal processing and communications," *Journal of the Franklin Institute*, special issue on *Modelling and Simulation in Advanced Communications*, vol. 347, no. 5, pp. 818–835, June 2010.
107. M. J. D. Powell, "Convergence properties of algorithms for nonlinear optimization," *SIAM Review*, vol. 28, no. 4, pp. 487–500, Dec. 1986.
108. S. S. Rangwala and D. A. Dornfeld, "Learning and optimization of machining operations using computing abilities of neural networks," *IEEE Transactions on Systems, Man and Cybernetics*, vol. 19, no. 2, pp. 299–314, Mar./Apr. 1989.
109. G. V. Reklaitis, A. Ravindran, and K. M. Ragsdell, *Engineering Optimization: Methods and Applications*. New York: Wiley-Interscience, 1983.
110. A. M. Rockett and J. C. Stevenson, "Karmarkar's algorithm: A method for solving large linear programming problems," *Byte*, vol. 12, no. 10, pp. 146–160, Sept. 1987.
111. H. L. Royden, *Real Analysis*, 3rd ed. New York: Macmillan Company, 1988.
112. W. Rudin, *Principles of Mathematical Analysis*, 3rd ed. New York: McGraw-Hill Book Co., 1976.
113. D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, vol. 1, *Foundations*. Cambridge, MA: MIT Press, 1986.
114. D. Russell, *Optimization Theory*. New York: W. A. Benjamin, 1970.
115. S. L. Salas and E. Hille, *Calculus: One and Several Variables*, 4th ed. New York: Wiley, 1982.

116. L. L. Scharf, L. T. McWhorter, E. K. P. Chong, J. S. Goldstein, and M. D. Zoltowski, "Algebraic equivalence of conjugate direction and multistage Wiener filters," in *Proceedings of the Eleventh Annual Workshop on Adaptive Sensor Array Processing (ASAP)*, Lexington, Massachusetts, Mar. 11–13, 2003.
117. L. L. Scharf, E. K. P. Chong, and Z. Zhang, "Algebraic equivalence of matrix conjugate direction and matrix multistage filters for estimating random vectors," in *Proceedings of the 43rd IEEE Conference on Decision and Control (CDC'04)*, Atlantis Resort, Paradise Island, Bahamas, Dec. 14–17, 2004, pp. 4175–4179.
118. L. L. Scharf, E. K. P. Chong, M. D. Zoltowski, J. S. Goldstein, and I. S. Reed, "Subspace expansion and the equivalence of conjugate direction and multistage Wiener filters," *IEEE Transactions on Signal Processing*, vol. 56, no. 10, pp. 5013–5019, Oct. 2008.
119. A. Schrijver, *Theory of Linear and Integer Programming*. New York: Wiley, 1986.
120. R. T. Seeley, *Calculus of Several Variables: An Introduction*. Glenview, IL: Scott, Foresman and Co., 1970.
121. J. R. Silvester, "Determinants of block matrices," *The Mathematical Gazette*, vol. 48, no. 51, pp. 460–467, Nov. 2000.
122. W. Spendley, G. R. Hext, and F. R. Himsorth, "Sequential application of simplex designs in optimization and evolutionary operation," *Technometrics*, vol. 4, pp. 441–461, 1962.
123. W. A. Spivey, *Linear Programming: An Introduction*. New York: Macmillan Company, 1963.
124. R. E. Stone and C. A. Tovey, "The simplex and projective scaling algorithms as iteratively reweighted least squares methods," *SIAM Review*, vol. 33, no. 2, pp. 220–237, June 1991.
125. G. Strang, *Introduction to Applied Mathematics*. Wellesley, MA: Wellesley-Cambridge Press, 1986.
126. G. Strang, *Linear Algebra and Its Applications*. New York: Academic Press, 1980.
127. T. W. Then and E. K. P. Chong, "Genetic algorithms in noisy environments," in *Proceedings of the 9th IEEE Symposium on Intelligent Control*, pp. 225–230, Aug. 1994.
128. L. Vandenberghe and S. Boyd, "Semidefinite programming," *SIAM Review*, vol. 38, no. 1, pp. 49–95, Mar. 1996.
129. P. P. Varaiya, *Notes on Optimization*. New York: Van Nostrand Reinhold Co., 1972.
130. D. J. Velleman, *How To Prove It: A Structured Approach*. Cambridge, England: Cambridge University Press, 1994.
131. L.-X. Wang, *A Course in Fuzzy Systems and Control*. Upper Saddle River, NJ: Prentice Hall, 1999.

132. B. Widrow and M. A. Lehr, "30 years of adaptive neural networks: Perceptron, madaline, and backpropagation," *Proceedings of the IEEE*, vol. 78, no. 9, pp. 1415–1442, Sept. 1990.
133. D. J. Wilde, *Optimum Seeking Methods*. Englewood Cliffs, NJ: Prentice Hall, 1964.
134. R. E. Williamson and H. F. Trotter, *Multivariable Mathematics, 2nd ed.* Englewood Cliffs, NJ: Prentice Hall, 1979.
135. W. I. Zangwill, *Nonlinear Programming: A Unified Approach*. Englewood Cliffs, NJ: Prentice Hall, 1969.
136. G. Zoutendijk, *Mathematical Programming Methods*. Amsterdam, The Netherlands: North-Holland, 1976.
137. J. M. Zurada, *Introduction to Artificial Neural Systems*. St. Paul, MN: West Publishing Co., 1992.



# Index

---

- Absolute value, 19  
Absolute value penalty function, 561  
Activation function, 254  
Active constraint, 483  
Adaline, 257, 271  
Adaptive linear element, 257  
Additivity, 19  
Affine function, 62, 63, 247, 524  
Affine matrix inequality, 524  
Affine scaling, 406, 407  
Affine scaling method, 406–411  
    artificial problem, 410  
    stopping criterion, 409  
    strictly interior feasible point, 406  
Algebraic Riccati inequality, 544  
Algorithm  
    affine scaling, 406–411  
    backpropagation, 253, 258–269  
    BFGS, 207–211  
    Broyden-Fletcher-Goldfarb-Shanno, *see* BFGS algorithm  
    complexity of, 306, 402  
    conjugate gradient, *see* Conjugate gradient algorithm  
    convergence of, *see* Convergence  
    Davidon-Fletcher-Powell, *see* DFP algorithm  
    DFP, 202–207  
    ellipsoid, *see* Khachiyan's method  
    exponential complexity, 402  
    fixed step size, 146, 151, 261  
    for constrained optimization, 545–567  
    genetic, 285–298  
    globally monotone, 157  
    gradient, 131–153  
    Gram-Schmidt, 177, 189  
    interior-point, 307, 403, 406, 411, 423  
    iterative, 124, 159. *See also* Search methods  
    Kaczmarz's, 232–236, 257  
    Karmarkar's, *see* Karmarkar's method  
    Khachiyan's, 306, 402–405, 418  
    Lagrangian, 553–560  
    naive random search, 279  
    Nelder-Mead, 274–278  
    particle swarm optimization, 282–285  
    polynomial complexity, 402

- probabilistic search, 278
- projected, 546, 556
- projected gradient, 407, 547, 549–553
- projected steepest descent, 551
- quasi-Newton, *see* Quasi-Newton methods
- randomized search, 278
- rank one, 197–202
- rank two, 202
- RLS, 227–232, 250
- secant method, 120–122, 172
- simplex, *see* Simplex method
- simulated annealing, 278–282
- single-rank symmetric, 197
- SRS, 197
- steepest descent, 133–141
- symmetric Huang family, 215
- variable metric, 202
- Widrow-Hoff, 257, 271
- zero finding, 118, 155
- Allocation**, 541
- Alphabet** in genetic algorithm, 286
- Argmin**, 82
- Armijo** backtracking algorithm, 126
- Armijo** condition, 125
- Armijo-Goldstein** condition, 125
- Artificial neural networks**, *see* Feedforward neural networks
- Artificial problem**
  - in affine scaling method, 410
  - in Karmarkar's method, 418
  - in simplex method, 361
- Associative**, 8, 9
- Asymmetric duality**, 379
- Augmented matrix**, 325
- Backpropagation algorithm**, 253, 258–269
  - as a gradient algorithm, 261
  - forward pass, 263
  - reverse pass, 263
- Ball**, 50
- Banana (Rosenbrock's) function**, 68, 159, 173, 191, 215
- Basic columns**, 325
- Basic feasible solution**, 305, 325, 327, 331, 349
- Basic solutions**, 324–327
- Basic variables**, 325, 347
- Basis**
  - definition of, 11
  - entering, 351
  - in linear equations, 325, 347, 348
- leaving**, 351
- natural**, 11
- orthogonal**, 29
- Beltrami**, 562
- Best-so-far**, 280, 283, 290
- BFGS algorithm**, 207–211
- Big-oh notation**, 74, 149
- Bisection method**, 116
- Bland's rule**, 360, 375
- Boltzmann**, 281
- Bolzano-Weierstrass theorem**, 58, 234
- Boundary**, 51
- Boundary point**, 51
- Bounded above**, 57
- Bounded below**, 57
- Bounded sequence**, 56, 57
- Bounded set**, 51
- Box constraint**, 546
- Bracketing**, 123, 129
- Brent's method**, 122
- Broyden**, 207
- Broyden-Fletcher-Goldfarb-Shanno algorithm**, *see* BFGS algorithm
- Canonical augmented matrix**, 346–348
- Canonical form**, 346
- Canonical representation**, 347
- Canonical representation of LMI**, 529
- Canonical tableau**, 358
- Carrier of polyhedron**, 53
- Cauchy-Schwarz inequality**, 20, 132, 207, 551
- Center of gravity**, 102, 275
- Centroid**, 102, 274
- Chain rule**, 67
- Characteristic equation**, 26
- Characteristic polynomial**, 26
- Chromosome in genetic algorithm**, 286
- Circuit**, 121, 253, 254, 311, 487
- Citation style**, 6
- Clairaut's theorem**, 65
- Closed set**, 51
- Column vector**, 7
- Combinatorial optimization**, 273, 282, 297
- Commutative**, 8
- Compact set**, 51
- Compatible matrix norm**, 36
- Complementarity**, 207
- Complementary slackness**, 390, 539
- Complex inner product**, 21
- Complex vector space**, 12
- Complexity of algorithm**, 306, 402
  - exponential, 402

- polynomial, 402
- Component of vector, 7
- Composite function, 67, 85, 90, 459
- Concave function, 512, *see Convex function*
- Condition number, 148
- Conjugate direction methods, 175–188
- Conjugate gradient algorithm
  - Fletcher-Reeves formula, 187
  - Hestenes-Stiefel formula, 187
  - nonquadratic problems, 186–188
  - Polak-Ribière formula, 187
  - Powell formula, 188
  - quadratic problems, 182–186
  - stopping criterion, 188
- Consistent linear inequalities, 404
- Constrained optimization, 82, 305, 449
- Constraint
  - active, 483
  - box, 546
  - convex, 517
  - equality, 450, 504
  - functional, 82
  - inactive, 483
  - inequality, 483, 561
  - set, 82
- Constraint set, 81. *See also Feasible set*
- Continuity, 21, 60, 453, 565
- Continuous function, 21, 60, 453, 565
- Continuously differentiable function, 65, 451, 513
- Contradiction, proof, 5
- Contraposition, proof, 4
- Contrapositive, 4
- Control system, 102, 472, 474, 475, 481, 553
- Convergence
  - fixed-step-size gradient algorithm, 146
  - globally convergent, 141
  - gradient algorithms, 143
  - Kaczmarz's algorithm, 233
  - linear, 148
  - locally convergent, 141, 554, 556
  - Newton's method, 165
  - of sequence of matrices, 59
  - order of, 148, 149, 152, 153, 157, 165
  - penalty method, 565
  - quadratic (second-order), 148
  - rate of, 141, 148
  - ratio, 148
  - steepest descent algorithm, 145
  - sublinear, 148
  - superlinear, 148
- Convergent sequence, 56
- Convex combination, 48, 297, 582
- Convex constraint, 517
- Convex function, 508–517
  - definition of, 509
  - differentiable, 513
  - equivalent definition of, 509
  - minimizers of, 517
  - optimization of, 517–536
  - quadratic, 512
  - strict, 512
  - twice differentiable, 515
- Convex optimization, 517–536
- Convex programming, *see Convex optimization*
- Convex set, 48–50
  - definition of, 48, 330, 509
  - extreme point, 50, 331, 332, 335
  - in definition of convex function, 509
  - polyhedron, 52, 317
  - polytope, 52, 316
  - properties of, 49
  - supporting hyperplane, 52, 316
- Cooling schedule, 281
- Coordinates, 11
- Cost function, 81
- Courant-Beltrami penalty function, 562, 570
- Cramer's rule, 33
- Crossing site, 287
- Crossover in genetic algorithm, 287
  - crossing site, 287
  - multiple-point crossover, 287, 300
  - one-point crossover, 287
- Cubic fit, 122
- Curve, 453
- Cutting-plane method, 435
- Cycling in simplex method, 360, 375
- Dantzig, 306
- Davidon, 202
- Davidon-Fletcher-Powell algorithm, *see DFP algorithm*
- Decision variable, 81, 305, 541
- Decomposition
  - direct sum, 30, 458
  - orthogonal, 30, 226
- Decreasing sequence, 55
- Degenerate basic feasible solution, 325, 328, 360
- DeMorgan's law, 3
- Derivative, 62, 83
  - partial, 64

- Derivative descent search, 154  
 Derivative matrix, 64  
 Descent property, 135, 141, 144, 167, 168, 193, 552  
 Determinant, 14  
 DFP algorithm, 202–207  
 Diagonal matrix, 27, 407, 421, 541, 571  
 Diet problem, 308, 381, 388  
 Differentiable curve, 453  
 Differentiable function, 62, 63, 453  
 Dimension, 11, 452  
 Direct sum decomposition, 30, 458  
 Directional derivative, 84  
 Discrete Fourier series, 225  
 Discrete-time linear system, 102, 248, 474, 481  
 Distributive, 9  
 Domination, 575  
 Dual linear program, 378, 403  
 Dual nonlinear program, 543  
 Dual quadratic program, 399  
 Duality
  - asymmetric, 379
  - dual nonlinear program, 543
  - dual problem, 378, 403
  - dual quadratic program, 399
  - dual vector, 378
  - duality theorem, 387, 543
  - in quasi-Newton methods, 207
  - Karush-Kuhn-Tucker conditions, 395, 539
  - linear programming, 377, 539
  - nonlinear programming, 542
  - primal nonlinear program, 542
  - primal problem, 378, 403
  - primal quadratic program, 399
  - quadratic programming, 399
  - symmetric, 378, 403
  - weak duality lemma, 385, 399, 543
- Duality theorem, 387, 543  
 Dyadic product, 197
- Eberhart, Russell, 282  
 Edge of polyhedron, 53  
 Eigenvalue
  - definition of, 26
  - maximal, 144
  - minimal, 144
  - of symmetric matrix, 28, 35
- Eigenvector
  - definition of, 26
  - of symmetric matrix, 28
  - orthogonal, 28
  - relation to  $\mathbf{Q}$ -conjugacy, 189
- Electric circuit, 311, 487  
 Elementary matrix
  - elementary row operation, 341
  - first kind, 340
  - second kind, 340
  - third kind, 341
- Elementary row operation, 341  
 Elitism in genetic algorithm, 290  
 Ellipsoid, 405  
 Ellipsoid algorithm, *see* Khachiyan's method  
 Encoding in genetic algorithm, 285, 290, 297  
 Entry of matrix, 12  
 Epigraph, 508  
 Equality constraint, 450, 504  
 Estimation, 121, 223, 244  
 Euclidean inner product, 19  
 Euclidean norm, 20  
 Evolution in genetic algorithm, 287  
 Exact penalty function, 566  
 Exclusive OR, *see* XOR  
 Expanding subspace theorem, 181  
 Exponential complexity, 402  
 Extreme point, 50, 331, 332, 334  
 Extremizer, 82
- Face of polyhedron, 53, 316  
 Farkas's transposition theorem, 396  
 Feasibility problem, 398, 526, 529  
 Feasible direction, 84, 519, 566  
 Feasible point, 305, 450  
 Feasible set, 81, 450  
 Feedforward neural networks, 253–269
  - activation function, 254
  - Adaline, 257
  - backpropagation algorithm, 253, 258–269
  - function approximation, 255
  - hidden layer, 254
  - input layer, 254
  - learning, 253
  - neuron, 253
  - output layer, 254
  - single-neuron training, 256–258
  - supervised learning, 256
  - training, 253
  - training set, 255
  - unsupervised learning, 256
  - weights, 253, 254
- Fibonacci method, 108–115  
 Fibonacci sequence, 109  
 First-order Lagrangian algorithm, 553  
 First-order necessary condition

- equality constraint (Lagrange), 460, 461
  - in convex optimization, 518
  - inequality constraint (KKT), 484
  - interior case, 86
  - set constraint, 85
- Fitness in genetic algorithm, 286
- Fitting straight line, 101, 221, 227, 246, 247
- Fixed point, 553
- Fixed step size, 146, 151, 261
- Fletcher, 187, 202, 207
- Fletcher-Reeves formula, 187
- Floor, 435
- FONC, *see* First-order necessary condition
- Fourier series, 225
- Frobenius norm, 36
- Full-rank factorization, 236
- Function
  - affine, 62, 63, 247, 524
  - banana, 68, 159, 173, 191, 215
  - composite, 67, 85, 90, 459
  - concave, 512, *see* Convex function
  - continuous, 21, 60, 453, 565
  - continuously differentiable, 65, 451, 513
  - convex, 508–517
  - cost, 81
  - derivative matrix of, 64
  - derivative of, 62, 83
  - differentiable, 62, 63, 453
  - directional derivative of, 84
  - gradient of, 65, 71, 83, 131
  - graph of, 71, 508
  - Jacobian matrix of, 64
  - Lagrangian, 464, 468, 543, 553
  - linear, *see* Linear transformation
  - matrix-valued, 60
  - maximum rate of decrease, 132
  - maximum rate of increase, 71, 131
  - notation, 5
  - objective, 81
  - partial derivative of, 64
  - penalty, 560
  - Powell, 162
  - Rosenbrock's, 68, 159, 173, 191, 215
  - sigmoid, 266
  - twice continuously differentiable, 65, 515
  - twice differentiable, 65, 454
  - uniformly continuous, 23
  - unimodal, 104
  - utility, 542
- Function approximation, 255
- Functional constraint, 82
- Fundamental theorem of algebra, 26
- Fundamental theorem of linear algebra, 41
- Fundamental theorem of LP, 328
- Fuzzy linear programming, 585
- Gale's transposition theorem, 397
- Gauss-Newton method, 171
- Generalized eigenvalue, 468, 534
- Generalized inverse, 236, 238
- Genetic algorithm, 285–298
  - alphabet, 286
  - analysis of, 291–297
  - best-so-far chromosome, 290
  - chromosome, 286
  - crossover, 287
  - elitism, 290
  - encoding, 285, 290, 297
  - evolution, 287
  - fitness, 286
  - initial population, 285
  - length of schema, 294
  - mating pool, 286
  - mutation, 288
  - offspring, 287
  - order of schema, 294
  - parents, 287
  - population size, 286
  - real-number, 297–298
  - representation scheme, 286
  - roulette-wheel scheme, 286
  - schema, 292
  - selection, 286
  - stopping criterion, 290
  - tournament scheme, 286
- Gibbs, 281
- Global minimizer, 82, 94, 473, 517–519, 521, 537, 552
- Globally convergent, 141
- Globally monotone algorithm, 157
- Golden section, 106
- Golden section search, 104–108
- Goldfarb, 207
- Goldstein condition, 125
- Gomory cut, 435
- Gomory cutting-plane method, 435
- Gordan's transposition theorem, 397
- Gradient, 65, 71, 83, 131
- Gradient descent algorithm, *see* Algorithm, gradient
- Gradient methods, 131–153

- backpropagation algorithm, 253, 258–269  
*constrained optimization, see Projected gradient method*  
 convergence of, 141–147  
 convergence rate of, 147–153  
 descent property, 135, 141, 144  
*equality constraints, see Lagrangian algorithms*  
 fixed step size, 145  
*inequality constraints, see Lagrangian algorithms*  
 Lagrangian, 553–560  
 order of convergence, 152  
 projected, 547, 549–553  
 stopping criterion, 135
- Gram matrix, 220  
 Gram-Schmidt, 177, 189  
 Grammian, 220  
 Graph, 71, 508  
 Greatest lower bound, 57
- Hačijan, *see* Khachiyan  
 Hadamard product, 283  
 Hajek, 282  
 Half-space, 46, 316  
   negative, 46  
   positive, 46  
 Hessian, 65, 468, 515  
 Hessian matrix, 83  
 Hestenes, Magnus, 187, 188  
 Hestenes-Stiefel formula, 187  
 Hidden layer in neural network, 254  
 Hoff, 257  
 Holland, John, 285  
 Homogeneity, 19, 21  
 Huang family, 215  
 Hyperplane  
   definition of, 46  
   supporting, 52, 316  
   tangent to graph, 71
- Identity matrix, 16  
 ILP, *see* Integer linear programming  
 Image of matrix, *see* Range of matrix  
 Implicit function theorem, 457  
 Impulse response, 248  
 Inactive constraint, 483  
 Inconsistent system of equations, 217  
 Increasing sequence, 55  
 Indefinite matrix, 35  
 Induced matrix norm, 36, 480  
 Induction, principle of, 5  
 Inequality constraint, 483, 561
- Infimum, *see* Greatest lower bound  
 Inner product  
   complex, 21  
   Euclidean, 19  
   properties of, 19  
 Innovation, 229  
 Input layer in neural network, 254  
 Integer linear programming, 427–444  
*Integer programming, see Integer linear programming*  
 Interior, 51  
 Interior point, 51  
 Interior-point method, 307, 403, 406, 411, 423  
 Inverse  
   continuity of, 60  
   matrix, 16  
 Inverse Hessian, 194  
 Inverse parabolic interpolation, 122  
 Invertible matrix, *see* Nonsingular matrix  
 Iterative algorithm, 124, *see* Search methods, 159
- Jacobian matrix, 64  
 Jordan form, 59
- Kaczmarz's algorithm, 232–236, 257  
 Kantorovich, 306  
 Karmarkar, 306, 403  
 Karmarkar's method, 306, 403, 411–423  
   artificial problem, 418  
   complexity, 403  
   Karmarkar's canonical form, 411–413, 415  
   Karmarkar's restricted problem, 414–415  
   projective transformation, 416, 424  
   simplex, 412  
   stopping criterion, 415, 420  
   strictly interior feasible point, 407, 416
- Karush-Kuhn-Tucker condition, *see* KKT condition  
 Karush-Kuhn-Tucker multiplier, *see* KKT multiplier  
 Karush-Kuhn-Tucker theorem, 484  
 Kennedy, James, 282  
 Kernel of matrix, *see* Nullspace of matrix  
 Khachiyan, 306, 402  
 Khachiyan's method, 306, 402–405, 418  
 KKT condition, 484, 485, 489, 521, 556, 594  
 KKT multiplier, 484, 490

- KKT theorem, 484
- Klee-Minty problem, 401
- Koopmans, 306
- Krylov subspace, 188
- Kuhn-Tucker condition, *see* KKT condition
- Lagrange condition, 460, 464, 520, 550, 553, 593
- Lagrange multiplier, 460, 463
- Lagrange's theorem, 460, 461
- Lagrangian algorithms, 553–560
- Lagrangian function, 464, 468, 543, 553
- Lanczos, Cornelius, 188
- Leading principal minor, 31
- Learning in neural network, 253
- Least squares, 217–227, 238
  - nonlinear, 169
- Least upper bound, 57
- Left pseudoinverse, 238
- Level set, 68, 131, 134
- Levenberg-Marquardt algorithm, 171
- Levenberg-Marquardt modification, 168
- Limit of sequence, 55
- Line fitting, 101, 221, 227, 246, 247
- Line search, 103, 124, 133, 167, 186, 188, 194, 209
- Line segment, 45, 48
- Linear combination, 10
- Linear convergence, 148
- Linear dynamical system, *see* Discrete-time linear system
- Linear equations
  - augmented matrix, 325
  - basic solution, 325
  - basis, 325, 347, 348
  - canonical augmented matrix, 347
  - canonical form, 346
  - degenerate basic solutions, 325
  - existence of solution, 17
  - inconsistent, 217
  - Kaczmarz's algorithm, 232–236
  - least-squares solution, 217, 218, 221
  - minimum-norm solution, 231, 241, 257, 473
  - overdetermined, 217
  - particular solution, 346
  - pivot, 349, 352, 364
  - solving in general, 217–243
  - solving using row operations, 339–346
- Linear function, *see* Linear transformation
- Linear inequalities
- consistent, 404
- in linear programming, 305, 307, 316
- Linear least squares, 217–227, 238
- Linear matrix inequality, 524, 541
- Linear programming
  - affine scaling method, 406–411
  - artificial problem in affine scaling method, 410
  - artificial problem in Karmarkar's method, 418
  - artificial problem in simplex method, 361
  - artificial variables in simplex method, 361
  - as constrained problem, 450
  - asymmetric duality, 379
  - basic columns, 325
  - basic feasible solution, 305, 325, 327, 331, 349
  - basic solutions, 324–327
  - basic variables, 325, 347
  - Bland's rule, 360, 375
  - brief history of LP, 305
  - canonical augmented matrix, 347
  - canonical tableau, 358
  - complementary slackness, 390, 539
  - cycling, 360, 375
  - degenerate basic feasible solution, 325, 328, 360
  - dual problem, 378, 403
  - duality, *see* Duality
  - duality theorem, 387
  - examples of, 100, 307–314
  - extreme point, 331, 332, 334
  - feasible solution, 325
  - fundamental theorem of LP, 328
  - fuzzy, 585
  - geometric view of, 330
  - integer linear programming, 427–444
  - interior-point method, 307, 403, 406, 411, 423
  - Karmarkar's method, *see* Karmarkar's method
  - Karush-Kuhn-Tucker condition, 395, 501, 539
  - Khachiyan's method, 306, 403–405, 418
  - Klee-Minty problem, 401
  - optimal basic feasible solution, 328
  - optimal feasible solution, 328
  - primal problem, 378, 403

- reduced cost coefficient, 353, 357, 358, 391
- revised simplex method, 364–368
- sensitivity, 398
- simplex method, 306, 339–368
- slack variable, 319
- standard form, 318, 324
- surplus variable, 319
- symmetric duality, 378, 403
- tableau, 358
- two-dimensional, 314
- two-phase affine scaling method, 409
- two-phase simplex method, 360–364
- uncertain, 584–591
- weak duality lemma, 385, 399
- Linear quadratic regulator, 474
- Linear regression, *see* Line fitting
- Linear space, *see* Vector space
- Linear transformation, 25, 63
- Linear variety, 47
- Linear-fractional LMIs, 534
- Linearly dependent, 10
- Linearly independent, 9, 176, 324, 451, 484
- Little-oh notation, 74, 85
- LMI, *see* Linear matrix inequality
- LMI solvers, 529
- LMI toolbox for MATLAB, 529, 536
- LMITOOL, 536
- Local minimizer, 82, 83, 85, 90, 92, 461, 469, 471, 484, 494, 517
- Locally convergent, 141, 554, 556
- Location parameter, 267
- Lower bound, 57
- LP, *see* Linear programming
- LQR, 474
- Lyapunov inequality, 527, 543
- MacDuffee, 241
- Markov chain, 397
- Mating pool in genetic algorithm, 286
- MATLAB, xiii, 127, 128, 136, 139, 159, 191, 215, 271, 290, 291, 299, 301, 337, 376, 424, 543
- LMI toolbox, 529, 536
- Matrix
  - affine matrix inequality, 524
  - compatible norm, 36
  - condition number, 148
  - continuous, 60
  - convergence of sequence, 59
  - definition of, 12
  - derivative, 64
  - determinant, 14
  - diagonal, 27, 407, 421, 541, 571
  - eigenvalue of, *see* Eigenvalue
  - eigenvector of, *see* Eigenvector
  - elementary, *see* Elementary matrix
  - entry of, 12
  - full-rank factorization, 236
  - function, matrix-valued, 60
  - game theory, 314
  - generalized inverse, 236, 238
  - Gram, 220
  - Hadamard product, 283
  - Hessian, 65, 83, 468, 515
  - identity, 16
  - image of, *see* Range of matrix
  - indefinite, 35
  - induced norm, 36, 480
  - inverse, 16
  - invertible, *see* Nonsingular matrix
  - Jacobian, 64
  - Jordan form, 59
  - kernel of, *see* Nullspace of matrix
  - leading principal minor of, 31
  - left pseudoinverse, 238
  - linear matrix inequality, 524, 541
  - minor of, 15, 428
  - Moore-Penrose inverse, 236, 237
  - negative definite, 35
  - negative semidefinite, 35
  - nonsingular, 16, 208, 218, 229, 325, 342
  - notation, 12
  - nullspace of, 30, 41, 227, 372, 406, 454
  - orthogonal, 29, 571
  - orthogonal projector, 30, 226, 406, 408, 423, 549
  - Penrose generalized inverse, 243
  - positive definite, 35
  - positive semidefinite, 35
  - principal minor of, 31
  - pseudoinverse, 236, 238
  - range of, 30, 41
  - rank of, 13–16
  - representation of linear transformation, 25
  - right pseudoinverse, 238
  - Schur complement, 526
  - Schur product, 283
  - sequence of, 59
  - series of, 60
  - similar, 26
  - square, 14

- stochastic, 397
- submatrix of, 432
- Sylvester's criterion, 31
- symmetric, 28, 35, 139
- totally unimodular, 432
- trace, 534
- transformation, 26
- transpose of, 12
- unimodular, 428
- Matrix norm**, 35–39
- Matrix-valued function**, 60
- Max**, 21
- Maximizer**, 82
- Mean value theorem**, 76, 554, 558
- MILP**, *see* Mixed integer linear programming
- Min**, 15, 82
- Minimax**, 582, 587, 592
- Minimizer**
  - description of, 81
  - global, 82, 94, 473, 517–519, 521, 537, 552
  - local, 82, 83, 85, 90, 461, 469, 484, 494, 517
  - Pareto, 575
  - strict global, 82
  - strict local, 82, 92, 102, 471, 494
- Minimum norm**, 231, 242, 257, 584, 592
- Minor**
  - definition of, 15, 428
  - leading principal, 31
  - principal, 31
- Minty**, 401
- Mixed integer linear programming**, 444
- Monotone sequence**, 55, 57
- Moore-Penrose inverse**, 236, 237
- Morrison**, 208, 229
- Multicriteria optimization**, 573
- Multiobjective optimization**, 573, 586
- Mutation** in genetic algorithm, 288
- Naive random search**, 279
- Natural basis**, 11
- Negative definite**
  - matrix, 35
  - quadratic form, 31
- Negative half-space**, 46
- Negative semidefinite**
  - matrix, 35
  - quadratic form, 31
- Neighborhood**, 50
- Nelder-Mead algorithm**, 274–278
  - centroid, 274
  - contraction, 276
- expansion**, 275
- Neural networks**, *see* Feedforward neural networks
- Neuron**, 253
- Newton's method**
  - convergence of, 165
  - descent direction, 167
  - descent property, 167
  - for nonlinear least squares, 168–171
  - Gauss-Newton method, 171
  - general, 161–171
  - Levenberg-Marquardt modification of, 168
  - modification of, 167
  - of tangents, 119
  - one-dimensional, 116–119
  - order of convergence, 165
- Newton-Raphson method**, *see* Newton's method
- Non-strict inequality**, 528
- Nondecreasing sequence**, 55
- Nondifferentiable optimization**, 523
- Nondifferentiable penalty function**, 566
- Nonincreasing sequence**, 55
- Nonlinear least squares**, 169
- Nonsingular matrix**, 16, 208, 218, 229, 325, 342
- Norm**
  - compatible, 36
  - Euclidean, 20
  - Frobenius, 36
  - general vector norm, 21
  - induced, 36, 480
  - matrix, 35–39
  - $p$ -norm, 21, 584
  - properties of, 21
- Normal**, 47, 70
- Normal plane**, 458
- Normal space**, 457, 458
- Notation**, 5
- Nullspace of matrix**, 30, 41, 227, 372, 406, 454
- Objective function**, 81
- Offspring** in genetic algorithm, 287
- One-dimensional search methods**, 103–126
- Open set**, 51
- Optimal basic feasible solution**, 328
- Optimal control**, 472, 475, 481, 482, 553
- Optimal feasible solution in LP**, 328
- Optimization**
  - combinatorial, 273, 282, 297
  - constrained, 82, 305, 449

- convex, 517–536
- linear, *see* Linear programming
- multicriteria, 573
- multiobjective, 573, 586
- nondifferentiable, 523
- semidefinite, 523
- unconstrained, *see* Unconstrained optimization
- vector, 573
- with equality constraints, 449, 553
- with inequality constraints, 483, 556
- with set constraint, 82, 502
- Optimization algorithm, *see* Search methods
- Order of convergence, 148, 149, 152, 153, 157, 165
- Order symbol, 73, 149
- Orthant, 416
- Orthogonal, 70
- Orthogonal basis, 29
- Orthogonal complement, 29, 226, 458
- Orthogonal decomposition, 30, 226
- Orthogonal matrix, 29, 571
- Orthogonal projection, 30, 219, 406, 408, 423
- Orthogonal projector, 30, 226, 406, 408, 423, 549
- Orthogonal vectors, 20
- Outer product, 197
- Output layer in neural network, 254
- Overdetermined system of equations, 217
- Parents in genetic algorithm, 287
- Pareto front, 575
- Pareto minimizer, 575
- Partial derivative, 64
- Particle swarm optimization, 282–285
- Particular solution, 346
- Penalty function, 560
- Penalty method, 560–567
  - absolute value penalty function, 561
  - convergence, 565
  - Courant-Beltrami penalty function, 562, 570
  - exact penalty function, 566
  - nondifferentiable penalty function, 566
  - penalty function, 560
  - penalty parameter, 560
- Penalty parameter, 560
- Penrose, *see* Moore-Penrose inverse
- Penrose generalized inverse, 243
- Perp, *see* Orthogonal complement
- Pivot, 349, 352, 364
- Polak-Ribière formula, 187
- Polyhedron
  - carrier of, 53
  - definition of, 52
  - edge of, 53
  - face of, 53, 316
  - in linear programming, 316–318
  - vertex of, 53
- Polynomial, characteristic, 26
- Polynomial complexity, 402
- Polytope
  - definition of, 52
  - in linear programming, 316
- Population in genetic algorithm, 285, 286
- Positive definite
  - matrix, 35
  - quadratic form, 31
  - relation to eigenvalues, 35
  - Sylvester's criterion, 31
- Positive half-space, 46
- Positive orthant, 416
- Positive semidefinite
  - matrix, 35
  - quadratic form, 31
  - relation to eigenvalues, 35
  - relation to principal minors, 35
- Positivity, 19, 21
- Powell, 162, 188, 202
- Powell formula, 188
- Powell function, 162
- Primal linear program, 378, 403
- Primal nonlinear program, 542
- Primal quadratic program, 399
- Primal-dual method, 377
- Principal minor, 31
- Principle of induction, 5
- Probabilistic search, 278
- Probability vector, 397, 540
- Product
  - dyadic, 197
  - inner, 19, 21
  - outer, 197
- Product rule, 68
- Projected algorithm, 546, 556
- Projected gradient method, 407, 547, 549–553
  - stopping criterion, 552
- Projected steepest descent algorithm, 551
- Projection, 297, *see* Orthogonal projection, 546
- Projective transformation, 416, 424

- Proof**
- contradiction (*reductio ad absurdum*), 5
  - contraposition, 4
  - direct method, 4
  - methods of, 3–5
  - principle of induction, 5
- Proportional fairness**, 541
- Pseudoinverse**, 236, 238
- Pythagorean theorem**, 21
- $Q$ -conjugate**
- definition of, 176
  - linear independence, 176
  - relation to eigenvectors, 189
  - relation to orthogonality, 189
- Quadratic convergence**, 148
- Quadratic fit**, 122, 128
- Quadratic form**
- convex, 512
  - definition of, 31
  - maximizing, 467, 471
  - negative definite, 31
  - negative semidefinite, 31
  - positive definite, 31
  - positive semidefinite, 31, 35
  - Sylvester’s criterion, 31
- Quadratic programming**, 399, 472, 481, 500
- Quasi-Newton methods**, 193–211
- approximating inverse Hessian, 194
  - BFGS algorithm, 207–211
  - complementarity, 207
  - conjugate direction property, 196
  - descent property, 193
  - DFP algorithm, 202–207
  - duality, 207
  - rank one formula, 197–202
  - rank two update, 202
  - single-rank symmetric, 197
  - symmetric Huang family, 215
  - variable metric algorithm, 202
- Randomized search**, 278
- Range of matrix**, 30, 41
- Rank of matrix**, 13–16
- Rank one formula**, 197–202
- Rank two update**, 202
- Rate of convergence**, 141, 148
- Ratio of convergence**, 148
- Rayleigh’s inequality**, 38, 92, 144, 146, 152
- Real vector space**, 7
- Recursive least-squares**, *see* RLS algorithm
- Reduced cost coefficient**, 353, 357, 358, 391
- Reductio ad absurdum**, 5
- Reeves**, 187
- Regular point**, 451, 456, 461, 464, 484, 554, 556
- Relative cost coefficient**, *see* Reduced cost coefficient
- Representation scheme in genetic algorithm**, 286
- Revised simplex method**, 364–368
- Revised tableau**, 365
- Ribière**, 187
- Riccati inequality**, 544
- Right pseudoinverse**, 238
- RLS algorithm**, 227–232, 250
- Rosenbrock’s function**, 68, 159, 173, 191, 215
- Roulette-wheel scheme**, 286
- Row operations**, 339–346
- Row vector**, 7
- Scalar**, 9
- Scale parameter**, 267
- Scheme in genetic algorithm**, 292
- length of, 294
  - order of, 294
- Schmidt**, *see* Gram-Schmidt
- Schur complement**, 526
- Schur product**, 283
- Schwarz**, *see* Cauchy-Schwarz inequality
- Schwarz’s theorem**, 65
- Scilab Consortium**, 536
- Search direction**, 124, 128, 167, 168
- Search methods**
- bisection method, 116
  - conjugate direction methods, 175–188
  - conjugate gradient algorithm, 182–188
  - constrained optimization, 545–567
  - derivative descent search, 154
  - Fibonacci, 108–115
  - general algorithm, 211
  - genetic algorithm, 285–298
  - Golden section, 104–108
  - gradient methods, 131–153
  - Kaczmarz’s algorithm, 232–236, 257
  - Lagrangian, 553–560
  - line search, 103, 124, 133, 167, 186, 188, 194, 209

- naive random search, 279
- Nelder-Mead algorithm, 274–278
- neural network training, 255
- Newton's method, 116–119, 161–171
- Newton-Raphson method, *see* Newton's method
- one-dimensional, 103–126
- particle swarm optimization, 282–285
- penalty method, 560–567
- probabilistic, 278
- projected, 546, 556
- projected gradient methods, 547, 549–553
- quasi-Newton methods, 193–211
- randomized, 278
- secant method, 120–122, 172
- simulated annealing algorithm, 278–282
- steepest descent method, 133–141
- Secant method, 120–122, 172
- Second-order necessary condition
  - equality constraints, 469
  - inequality constraints, 494
  - interior case, 90
  - set constraint, 90
- Second-order sufficient condition
  - equality constraints, 470
  - inequality constraints, 494
  - interior case, 92
  - set constraint, 102
- Selection in genetic algorithm, 286
- Semidefinite programming, 523
- Sensitivity, 398
- Sequence
  - bounded, 56, 57
  - bounded above, 57
  - bounded below, 57
  - convergent, 56
  - decreasing, 55
  - Fibonacci, 109
  - greatest lower bound, 57
  - increasing, 55
  - least upper bound, 57
  - limit, 55
  - lower bound, 57
  - monotone, 55, 57
  - nondecreasing, 55, 564
  - nonincreasing, 55
  - of matrices, 59
  - of real numbers, 55
  - order of convergence, 148, 149, 152, 153, 157, 165
- subsequence of, 58
- upper bound, 57
- Set
  - boundary of, 51
  - bounded, 51
  - closed, 51
  - compact, 51
  - constraint, 81, *see* Feasible set
  - convex, *see* Convex set
  - feasible, 81, 450
  - interior of, 51
  - minus, 5
  - notation, 5
  - open, 51
  - simplex, 274, 412
  - subset of, 5
- Set constraint, 82, 502
- Shanno, 207
- Sherman-Morrison formula, 208, 229
- Sherman-Morrison-Woodbury formula, 229
- Shift parameter, 267
- Sigmoid, 266
- Signal-to-interference ratio, 88, 98
- Similar matrices, 26
- Simplex, 274, 412
- Simplex algorithm, *see* Simplex method
- Simplex method, 306, 339–368
  - algorithm, 349–356
  - artificial problem, 361
  - artificial variables, 361
  - Bland's rule, 360, 375
  - canonical augmented matrix, 346–348
  - canonical tableau, 358
  - complexity, 402
  - cycling, 360, 375
  - exponential complexity, 402
  - integer linear programming, 427–444
  - matrix form, 356–360
  - reduced cost coefficient, 353, 357, 358, 391
  - revised simplex method, 364–368
  - revised tableau, 365
  - row operations, 339–346
  - stopping criterion, 351, 372
  - tableau, 358
  - two-phase, 360–364
  - updating augmented matrix, 348–349
  - updating canonical tableau, 358
- Simulating annealing algorithm, 278–282

- Simultaneous equations, *see* Linear equations  
 Single-rank symmetric algorithm, 197  
 Singular value decomposition, 571  
 Slack variable, 319  
 SONC, *see* Second-order necessary condition  
 SOSC, *see* Second-order sufficient condition  
 Span, 10  
 Sphere, 425  
 Square matrix, 14  
 SRS algorithm, 197  
 Standard form linear program, 318, 324  
 Statement  
     biconditional, 4  
     conditional, 3  
 Steepest ascent, 71  
 Steepest ascent method, *see* Steepest descent method  
 Steepest descent  
     order of convergence, 152  
 Steepest descent method, 133–141  
     for constrained optimization, 551  
     for quadratic, 139  
     projected, 551  
 Step response, 248  
 Step size, 132, 146, 156, 168, 233, 261, 547, 551  
 Stiefel, Eduard, 187, 188  
 Stochastic matrix, 397  
 Stopping criterion  
     affine scaling method, 409  
     conjugate gradient method, 188  
     genetic algorithm, 290  
     gradient method, 135  
     Karmarkar's method, 415, 420  
     line search, 128  
     projected gradient method, 552  
     simplex method, 351, 372  
 Strict inequality, 528  
 Strictly interior feasible point, 406, 407, 416  
 Strong Wolfe condition, 126  
 Structured representation of LMI, 529  
 Subgradient, 515, 537  
 Sublinear convergence, 148  
 Submatrix, 432  
 Subsequence, 58  
 Subset, 5  
 Subspace, 10  
 Superlinear convergence, 148  
 Supervised learning, 256  
 Supporting hyperplane, 52, 316  
 Supremum, *see* Least upper bound  
 Surface, 451, 454  
 Surplus variable, 319  
 SVD, *see* Singular value decomposition  
 Sylvester's criterion, 31  
 Symmetric duality, 378, 403  
 Symmetric Huang family, 215  
 Symmetric matrix, 28, 35, 139  
 Symmetry, 19  
 Tableau in linear programming, 358  
 Tangent line, 70  
 Tangent plane, 70, 454  
 Tangent space, 454, 456, 458  
 Tangent vector, 65, 70, 454, 456  
 Taylor series, 72–76, 162, 165, 551. *See also* Taylor's theorem  
 Taylor's formula, 72, 75. *See also* Taylor's theorem  
 Taylor's theorem, 72, 85, 90, 92, 132, 515, 516  
 Temperature schedule, 281  
 Termination criterion, *see* Stopping criterion  
 Third-order necessary condition, 94  
 Third-order sufficient condition, 94  
 Threshold, 267  
 Totally unimodular, 432  
 Tournament scheme, 286  
 Trace, 534  
 Training of neural network, 253  
 Training set, 255  
 Transformation  
     affine scaling, 407  
     linear, 25, 63  
     matrix, 26  
     matrix representation of, 25  
     projective, 416, 424  
 Transportation problem, 306, 310  
 Transpose  
     matrix, 12  
     vector, 8  
 Transposition theorems, 396  
 Traveling salesperson problem, 282  
 Triangle inequality, 21  
 Truth table, 3  
 Tucker, *see* KKT condition  
 Twice continuously differentiable function, 65, 515  
 Twice differentiable function, 65, 454  
 Two-dimensional linear program, 314  
 Two-phase affine scaling method, 409  
 Two-phase simplex method, 360–364

- Uncertainty range, 106
- Unconstrained optimization
  - basics of, 82
  - conditions for, 83–92
- Uniform continuity, 23
- Uniformly continuous function, 23
- Unimodal, 104
- Unimodular, 428
- Unimodular, totally 432
- Unsupervised learning, 256
- Upper bound, 57
- Utility function, 542
  
- Variable metric algorithm, 202
- Variety, linear, 47
- Vector
  - column, 7
  - complex, 12
  - component of, 7
  - convex combination, 48, 297, 582
  - definition of, 7
  - difference, 8
  - field, 65
  - linear combination, 10
  - linearly dependent, 10
  - linearly independent, 9, 176, 324, 451, 484
  - normal, 47
  - orthogonal, 20
  - probability, 397, 540
  - row, 7
  - tangent, 65, 70, 454, 456
  - transpose of, 8
  - zero vector, 8
- Vector field, 65
- Vector optimization, 573
- Vector space
  - basis for, 11
  - complex, 12
  - definition of, 7
  - dimension of, 11
  - real, 7
  - subspace of, 10
- Vertex, 53, 274, 406
  
- Weak duality lemma, 385, 399, 543
- Weierstrass theorem, 36, 51
- Weighted sum, 474, 582, 592
- Weights in neural network, 253, 254
- Widrow, 257
- Widrow-Hoff algorithm, 257, 271
- Wiener filter, 188, 245
- Wolfe condition, 125
- Woodbury, 229
  
- XOR, 268, 271
  
- YALMIP, 536
- Yet Another LMI Package (YALMIP), 536
  
- Zero finding, 118, 155
- Zero matrix, 22, 30, 36
- Zero vector, 8

## WILEY SERIES IN DISCRETE MATHEMATICS AND OPTIMIZATION

- AARTS AND KORST • Simulated Annealing and Boltzmann Machines: A Stochastic Approach to Combinatorial Optimization and Neural Computing
- AARTS AND LENSTRA • Local Search in Combinatorial Optimization
- ALON AND SPENCER • The Probabilistic Method, Third Edition
- ANDERSON AND NASH • Linear Programming in Infinite-Dimensional Spaces: Theory and Application
- ARLINGHAUS, ARLINGHAUS, AND HARARY • Graph Theory and Geography: An Interactive View E-Book
- AZENCOTT • Simulated Annealing: Parallelization Techniques
- BARTHÉLEMY AND GUÉNOCHE • Trees and Proximity Representations
- BAZARRA, JARVIS, AND SHERALI • Linear Programming and Network Flows
- BRUEN AND FORCINITO • Cryptography, Information Theory, and Error-Correction: A Handbook for the 21st Century
- CHANDRU AND HOOKER • Optimization Methods for Logical Inference
- CHONG AND ŽAK • An Introduction to Optimization, Fourth Edition
- COFFMAN AND LUEKER • Probabilistic Analysis of Packing and Partitioning Algorithms
- COOK, CUNNINGHAM, PULLEYBLANK, AND SCHRIJVER • Combinatorial Optimization
- DASKIN • Network and Discrete Location: Modes, Algorithms and Applications
- DINITZ AND STINSON • Contemporary Design Theory: A Collection of Surveys
- DU AND KO • Theory of Computational Complexity
- ERICKSON • Introduction to Combinatorics
- GLOVER, KLINGHAM, AND PHILLIPS • Network Models in Optimization and Their Practical Problems
- GOLSHTEIN AND TRETYAKOV • Modified Lagrangians and Monotone Maps in Optimization
- GONDTRAN AND MINOUX • Graphs and Algorithms (*Translated by S. Vajdā*)
- GRAHAM, ROTHSCILD, AND SPENCER • Ramsey Theory, Second Edition
- GROSS AND TUCKER • Topological Graph Theory
- HALL • Combinatorial Theory, Second Edition
- HOOKER • Logic-Based Methods for Optimization: Combining Optimization and Constraint Satisfaction
- IMRICH AND KLAVŽAR • Product Graphs: Structure and Recognition
- JANSON, LUCZAK, AND RUCINSKI • Random Graphs
- JENSEN AND TOFT • Graph Coloring Problems
- KAPLAN • Maxima and Minima with Applications: Practical Optimization and Duality
- LAWLER, LENSTRA, RINNOOY KAN, AND SHMOYS, Editors • The Traveling Salesman Problem: A Guided Tour of Combinatorial Optimization
- LAYWINE AND MULLEN • Discrete Mathematics Using Latin Squares
- LEVITIN • Perturbation Theory in Mathematical Programming Applications
- MAHMOUD • Evolution of Random Search Trees
- MAHMOUD • Sorting: A Distribution Theory
- MARTELLI • Introduction to Discrete Dynamical Systems and Chaos
- MARTELLO AND TOTH • Knapsack Problems: Algorithms and Computer Implementations
- McALOON AND TRETКОFF • Optimization and Computational Logic
- MERRIS • Combinatorics, Second Edition
- MERRIS • Graph Theory
- MINC • Nonnegative Matrices
- MINOUX • Mathematical Programming: Theory and Algorithms (*Translated by S. Vajdā*)
- MIRCHANDANI AND FRANCIS, Editors • Discrete Location Theory
- NEMHAUSER AND WOLSEY • Integer and Combinatorial Optimization
- NEMIROVSKY AND YUDIN • Problem Complexity and Method Efficiency in Optimization  
*(Translated by E. R. Dawson)*

PACH AND AGARWAL • Combinatorial Geometry  
PLESS • Introduction to the Theory of Error-Correcting Codes, Third Edition  
ROOS AND VIAL • Ph. Theory and Algorithms for Linear Optimization: An Interior Point Approach  
SCHEINERMAN AND ULLMAN • Fractional Graph Theory: A Rational Approach to the Theory of  
Graphs  
SCHIFF • Cellular Automata: A Discrete View of the World  
SCHRIVVER • Theory of Linear and Integer Programming  
SPALL • Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control  
STIEBITZ, SCHEIDE, TOFT, AND FAVRHOLDT • Graph Edge Coloring: Vizing's Theorem and  
Goldberg's Conjecture  
SZPANKOWSKI • Average Case Analysis of Algorithms on Sequences  
TOMESCU • Problems in Combinatorics and Graph Theory (*Translated by R. A. Melter*)  
TUCKER • Applied Combinatorics, Second Edition  
WOLSEY • Integer Programming  
YE • Interior Point Algorithms: Theory and Analysis