

# Investment Science

DAVID G. LUENBERGER

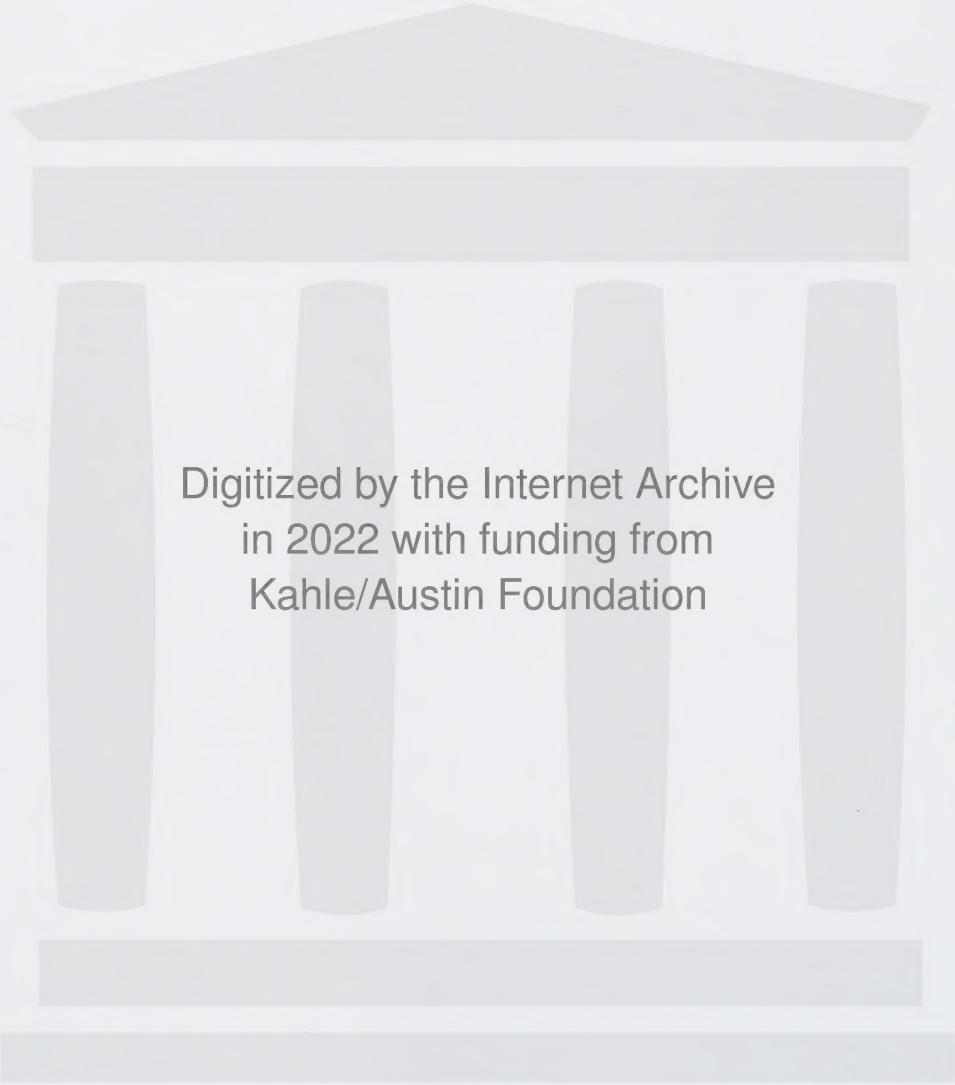
SECOND EDITION



OXFORD  
UNIVERSITY PRESS







Digitized by the Internet Archive  
in 2022 with funding from  
Kahle/Austin Foundation

[https://archive.org/details/investmentscienc0000luen\\_x4b7](https://archive.org/details/investmentscienc0000luen_x4b7)

# **INVESTMENT SCIENCE**

---



# INVESTMENT SCIENCE

SECOND EDITION

DAVID G. LUENBERGER

STANFORD UNIVERSITY

New York      Oxford  
OXFORD UNIVERSITY PRESS

Oxford University Press is a department of the University of Oxford.  
It furthers the University's objective of excellence in research, scholarship,  
and education by publishing worldwide.

Oxford New York

Auckland Cape Town Dar es Salaam Hong Kong Karachi  
Kuala Lumpur Madrid Melbourne Mexico City Nairobi  
New Delhi Shanghai Taipei Toronto

With offices in

Argentina Austria Brazil Chile Czech Republic France Greece  
Guatemala Hungary Italy Japan Poland Portugal Singapore  
South Korea Switzerland Thailand Turkey Ukraine Vietnam

Copyright © 2014 and 1998 by David G. Luenberger

For titles covered by Section 112 of the US Higher Education Opportunity Act, please visit [www.oup.com/us/he](http://www.oup.com/us/he) for the latest information about pricing and alternate formats.

Published by Oxford University Press.  
198 Madison Avenue, New York, NY 10016  
[www.oup.com](http://www.oup.com)

Oxford is a registered trademark of Oxford University Press

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

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

Luenberger, David G., 1937—

Investment science / David G. Luenberger. — Second Edition.

pages cm

ISBN 978-0-19-974008-6

1. Investments—Mathematical models. 2. Investment analysis—Mathematical models.
3. Cash flow—Mathematical models. 4. Interest rates—Mathematical models.
5. Derivative securities—Mathematical models. I. Title.

HG4515.2.L84 2013

332.6—dc23

2012047878

To My Family



# BRIEF CONTENTS

<b>PREFACE</b>	xxi
<b>Chapter 1 INTRODUCTION</b>	1
<b>Part I: DETERMINISTIC CASH FLOW STREAMS</b>	
<b>Chapter 2 THE BASIC THEORY OF INTEREST</b>	15
<b>Chapter 3 FIXED-INCOME SECURITIES</b>	42
<b>Chapter 4 THE TERM STRUCTURE OF INTEREST RATES</b>	76
<b>Chapter 5 APPLIED INTEREST RATE ANALYSIS</b>	107
<b>Part II: SINGLE-PERIOD RANDOM CASH FLOWS</b>	
<b>Chapter 6 MEAN-VARIANCE PORTFOLIO THEORY</b>	143
<b>Chapter 7 THE CAPITAL ASSET PRICING MODEL</b>	180
<b>Chapter 8 OTHER PRICING MODELS</b>	213
<b>Chapter 9 DATA AND STATISTICS</b>	235
<b>Chapter 10 RISK MEASURES</b>	257
<b>Chapter 11 GENERAL PRINCIPLES</b>	279
<b>Part III: DERIVATIVE SECURITIES</b>	
<b>Chapter 12 FORWARDS, FUTURES, AND SWAPS</b>	315
<b>Chapter 13 MODELS OF ASSET DYNAMICS</b>	350
<b>Chapter 14 BASIC OPTIONS THEORY</b>	374
<b>Chapter 15 ADDITIONAL OPTIONS TOPICS</b>	410
<b>Chapter 16 INTEREST RATE DERIVATIVES</b>	448
<b>Chapter 17 CREDIT RISK</b>	483
<b>Part IV: GENERAL CASH FLOW STREAMS</b>	
<b>Chapter 18 OPTIMAL PORTFOLIO GROWTH</b>	517
<b>Chapter 19 GENERAL INVESTMENT EVALUATION</b>	547

<b>Appendix A BASIC PROBABILITY THEORY</b>	579
<b>Appendix B CALCULUS AND OPTIMIZATION</b>	583
<b>ANSWERS TO EXERCISES</b>	588
<b>INDEX</b>	594

# CONTENTS

<b>PREFACE</b>	xxi
<b>Chapter 1 INTRODUCTION</b>	1
1.1 Cash Flows	2
1.2 Investments and Markets	3
The Comparison Principle	4
Arbitrage	4
Dynamics	5
Risk Aversion	5
1.3 Typical Investment Problems	6
Pricing	6
Hedging	7
Risk Assessment and Management	8
Pure Investment	8
Other Problems	9
1.4 Organization of the Book	9
Deterministic Cash Flow Streams	9
Single-Period Random Cash Flow Streams	10
Derivative Assets	10
General Cash Flow Streams	11
<b>Part I: DETERMINISTIC CASH FLOW STREAMS</b>	
<b>Chapter 2 THE BASIC THEORY OF INTEREST</b>	15
2.1 Principal and Interest	15
Simple Interest	15
Compound Interest	16
Compounding at Various Intervals	17
Continuous Compounding	18
Debt	19
Money Markets	19
2.2 Present Value	20
2.3 Present and Future Values of Streams	21
The Ideal Bank	21
Future Value	21
Present Value	22

	Frequent and Continuous Compounding	23
	Present Value and an Ideal Bank	23
2.4	Internal Rate of Return	24
2.5	Evaluation Criteria	26
	Net Present Value	27
	Internal Rate of Return	28
	Discussion of the Criteria	28
2.6	Applications and Extensions*	30
	Net Flows	30
	Cycle Problems	31
	Taxes	33
	Inflation	34
2.7	Summary	36
	Exercises	37
	References	41
<b>Chapter 3 FIXED-INCOME SECURITIES</b>		42
3.1	The Market for Future Cash	43
	Savings Deposits	43
	Money Market Instruments	44
	U.S. Government Securities	44
	Other Bonds	45
	Mortgages	46
	Annuities	46
3.2	Value Formulas	46
	Perpetual Annuities	47
	Finite-Life Streams	48
	Running Amortization*	50
	Annual Worth*	51
3.3	Bond Details	52
	Quality Ratings	53
3.4	Yield	54
	Qualitative Nature of Price–Yield Curves	55
	Other Yield Measures	58
3.5	Duration	59
	Interest Duration	60
	Macaulay Duration	60
	Explicit Formula*	61
	Qualitative Properties of Duration*	61
	Duration and Sensitivity	62
	Duration of a Portfolio	64
3.6	Immunization	65
3.7	Convexity*	68
3.8	Summary	69
	Exercises	71
	References	74

<b>Chapter 4 THE TERM STRUCTURE OF INTEREST RATES</b>	76
4.1 The Yield Curve	76
4.2 The Term Structure	78
Spot Rates	78
Discount Factors and Present Value	79
Determining the Spot Rate	81
4.3 Forward Rates	82
4.4 Term Structure Explanations	85
Expectations Theory	85
Liquidity Preference	86
Market Segmentation	87
Discussion	87
4.5 Expectations Dynamics	88
Spot Rate Forecasts	88
Discount Factors	89
Short Rates	90
Invariance Theorem	91
4.6 Running Present Value	92
4.7 Floating-Rate Bonds	95
4.8 Duration	96
Fisher-Weil Duration	96
Discrete-Time Compounding*	97
4.9 Immunization	98
4.10 Summary	100
Exercises	102
References	106
<b>Chapter 5 APPLIED INTEREST RATE ANALYSIS</b>	107
5.1 Capital Budgeting	108
Independent Projects	108
Interdependent Projects*	111
5.2 Optimal Portfolios	113
The Cash Matching Problem	114
5.3 Dynamic Cash Flow Processes	116
Representation of Dynamic Choice	117
Cash Flows in Graphs	119
5.4 Optimal Management	120
Running Dynamic Programming	120
Examples	123
5.5 The Harmony Theorem*	128
5.6 Valuation of a Firm*	130
Dividend Discount Models	130
Free Cash Flow*	132
5.7 Summary	134
Exercises	136
References	139

**Part II: SINGLE-PERIOD RANDOM CASH FLOWS**

<b>Chapter 6 MEAN-VARIANCE PORTFOLIO THEORY</b>	143
6.1 Asset Return	144
Short Sales	144
Portfolio Return	146
6.2 Random Variables	147
Expected Value	148
Variance	149
Several Random Variables	150
Covariance	150
Variance of a Sum	152
6.3 Random Returns	152
Mean–Standard Deviation Diagram	155
6.4 Portfolio Mean and Variance	156
Mean Return of a Portfolio	156
Variance of Portfolio Return	156
Diversification*	157
Diagram of a Portfolio	159
6.5 The Feasible Set	161
The Minimum-Variance Set and the Efficient Frontier	162
6.6 The Markowitz Model	164
Solution of the Markowitz Problem*	165
Nonnegativity Constraints*	168
6.7 The Two-Fund Theorem*	168
6.8 Inclusion of a Risk-Free Asset	171
6.9 The One-Fund Theorem	173
Solution Method*	173
Explicit Solution	175
6.10 Summary	175
Exercises	176
References	179
<b>Chapter 7 THE CAPITAL ASSET PRICING MODEL</b>	180
7.1 Market Equilibrium	180
7.2 The Capital Market Line	182
7.3 The Pricing Model	184
Betas of Common Stocks	187
Beta of a Portfolio	187
7.4 The Security Market Line	187
Systematic Risk	189
7.5 Investment Implications	190
7.6 Performance Evaluation	191
7.7 CAPM as a Pricing Formula	194
Linearity of Pricing and the Certainty Equivalent Form	196
7.8 Project Choice*	198

7.9	Projection Pricing	200
	Minimum Norm Pricing*	202
7.10	Correlation Pricing	203
7.11	Summary	206
	Exercises	207
	References	211
<b>Chapter 8 OTHER PRICING MODELS</b>		213
8.1	Introduction	213
8.2	Factor Models	213
	Single-Factor Model	214
	Portfolio Parameters	215
	Multifactor Models	219
	Selection of Factors	219
8.3	The CAPM as a Factor Model	220
	The Characteristic Line	221
8.4	Arbitrage Pricing Theory*	223
	Simple Version of APT	223
	Well-Diversified Portfolios	225
	General APT	226
	APT and CAPM	227
8.5	Projection Pricing with Factors	227
8.6	A Multiperiod Fallacy	229
8.7	Summary	230
	Exercises	232
	References	234
<b>Chapter 9 DATA AND STATISTICS</b>		235
9.1	Basic Estimation Methods	235
	Period-Length Effects	236
	Mean Blur	238
9.2	Estimation of Other Parameters	240
	Estimation of $\sigma$	240
	$\alpha$ Blur	241
9.3	The Effect of Estimation Errors	242
	Three Views	243
	Maximum Tangent	245
	Compounding Effect	248
9.4	Conservative Approaches	248
	Better Estimates*	249
9.5	Tilting Away From Equilibrium*	250
9.6	Summary	252
	Exercises	253
	References	255
<b>Chapter 10 RISK MEASURES</b>		257
10.1	Value at Risk	258

Properties of VaR	260
Capital Requirement	260
10.2 Computation of Value at Risk	261
Model-Based Method	261
Other Models	264
Shortcut for Discrete Distributions	264
Empirical Approach for Market Risk*	265
10.3 Criticisms of VaR	266
Diversification Failure	266
Poor Assessment of Risk	267
Discontinuous Value	268
10.4 Coherent Risk Measures	269
10.5 Conditional Value at Risk	270
10.6 Coherent Characterization*	272
10.7 Convexity*	274
10.8 Summary	275
Exercises	275
References	277
<b>Chapter 11 GENERAL PRINCIPLES</b>	279
11.1 Introduction	279
11.2 Utility Functions	279
Equivalent Utility Functions	281
11.3 Risk Aversion	282
Derivatives	284
Risk Aversion Coefficients	284
Certainty Equivalent	284
11.4 Specification of Utility Functions*	285
Direct Measurement of Utility	285
Parameter Families	287
Questionnaire Method	288
11.5 Utility Functions and the Mean–Variance Criterion*	288
Quadratic Utility	288
Normal Returns	290
11.6 Linear Pricing	291
Type A Arbitrage	291
Portfolios	292
Type B Arbitrage	292
11.7 Portfolio Choice	293
11.8 Arbitrage Bounds	296
11.9 Zero-Level Pricing	297
11.10 Log-Optimal Pricing*	299
11.11 Finite State Models	301
Completeness	302
State Prices	302
Positive State Prices	302

11.12 Risk-Neutral Pricing	304
11.13 Summary	306
Exercises	308
References	311

### **Part III: DERIVATIVE SECURITIES**

<b>Chapter 12 FORWARDS, FUTURES, AND SWAPS</b>	315
12.1 Pricing Principles	316
12.2 Forward Contracts	318
Forward Interest Rates	319
12.3 Forward Prices	319
Costs of Carry	322
Tight Markets	324
Investment Assets	325
12.4 The Value of a Forward Contract	326
12.5 Swaps*	327
Value of a Commodity Swap	327
Value of an Interest Rate Swap	329
12.6 Basics of Futures Contracts	329
12.7 Futures Prices	332
12.8 Relation to Expected Spot Price*	335
12.9 The Perfect Hedge	336
12.10 The Minimum-Variance Hedge	336
12.11 Optimal Hedging*	340
12.12 Hedging Nonlinear Risk*	341
12.13 Summary	345
Exercises	346
References	349
<b>Chapter 13 MODELS OF ASSET DYNAMICS</b>	350
13.1 Binomial Lattice Model	351
13.2 The Additive Model	353
Normal Price Distribution	354
13.3 The Multiplicative Model	355
Lognormal Prices	355
Real Stock Distributions	356
13.4 Typical Parameter Values*	357
13.5 Lognormal Random Variables	358
13.6 Random Walks and Wiener Processes	359
Generalized Wiener Processes and Ito Processes	361
13.7 A Stock Price Process	362
Lognormal Prices	363
Standard Ito Form	363
Simulation	365
13.8 Ito's Lemma*	366

13.9 Binomial Lattice Revisited	368
13.10 Summary	370
Exercises	370
References	373
<b>Chapter 14 BASIC OPTIONS THEORY</b>	374
14.1 Option Concepts	375
14.2 The Nature of Option Values	377
Time Value of Options	379
Other Factors Affecting the Value of Options	379
14.3 Option Combinations and Put–Call Parity	380
Put–Call Parity	381
14.4 Early Exercise	382
14.5 Single-Period Binomial Options Theory	383
14.6 Multiperiod Options	386
No Early Exercise*	389
14.7 More General Binomial Problems	389
Put Options	389
Dividend and Term Structure Problems*	391
Futures Options*	391
14.8 Evaluating Real Investment Opportunities	393
Real Options	397
Linear Pricing	399
14.9 General Risk-Neutral Pricing*	401
14.10 Three-principle Power	402
Decomposition of the Pricing Principles	403
14.11 Summary	403
Exercises	404
References	408
<b>Chapter 15 ADDITIONAL OPTIONS TOPICS</b>	410
15.1 Introduction	410
15.2 The Black–Scholes Equation	410
Proof of the Black–Scholes Equation*	412
Self-Financing Strategies*	414
15.3 Call Option Formula	414
15.4 Risk-Neutral Valuation*	416
15.5 Delta	417
15.6 Replication, Synthetic Options, and Portfolio Insurance*	419
15.7 Volatility Smiles	422
Equality of Implied Volatilities	423
Risk-Neutral Probability Density*	424
15.8 Computational Methods	425
Monte Carlo Simulation	426
Finite-Difference Methods	427
Binomial and Trinomial Lattices	429

15.9 Exotic Options	431
Pricing*	433
15.10 Comparison of Methods	434
15.11 Storage Costs and Dividends*	435
Binomial Form	435
Brownian Motion Form*	436
15.12 Martingale Pricing*	437
15.13 Axioms and Black–Scholes	438
Market Price of Risk	440
15.14 Summary	440
Exercises	442
References	446
<b>Chapter 16 INTEREST RATE DERIVATIVES</b>	448
16.1 Examples of Interest Rate Derivatives	448
16.2 The Need for a Theory	450
16.3 The Binomial Approach	451
Implied Term Structure	452
No Arbitrage Opportunities	454
16.4 Pricing Applications	455
Bond Derivatives	455
Forwards and Futures*	455
Futures*	457
16.5 Leveling and Adjustable-Rate Loans*	457
Adjustable-Rate Loans	458
16.6 The Forward Equation	461
16.7 Matching the Term Structure	464
The Ho–Lee Model	464
The Black–Derman–Toy Model	465
Matching Implied Volatilities	465
16.8 Immunization	467
16.9 Collateralized Mortgage Obligations*	469
16.10 Models of Interest Rate Dynamics*	473
16.11 Continuous-Time Solutions*	474
The Backward Equation	475
Affine Processes*	476
Risk-Neutral Pricing Formula	477
16.12 Extensions	477
16.13 Summary	478
Exercises	479
References	482
<b>Chapter 17 CREDIT RISK</b>	483
17.1 The Classic Merton Model	484
Probability of Default	486
Credit Spread	486

17.2	First Passage Times	487
	Lattice Methods	488
	Early Default*	490
	Coupons*	491
17.3	Rating Methods	492
17.4	Intensity (Reduced-Form) Model	493
	Poisson Processes	493
	Inhomogeneous Process	495
17.5	Stochastic Intensity Model*	495
17.6	Intermediate Receipts	496
17.7	Analytically Tractable Cox Processes	497
	Model Fitting	497
17.8	Simulation	498
	Direct Simulation	498
	A Better Way	499
17.9	Lattice Methods	500
17.10	Correlated Defaults	503
17.11	Credit Derivatives	505
	Bonds and Loans	506
	Credit Default Swaps (CDS's)	506
	Forwards and Options on CDS's	508
	Total Return Swaps (TRS's)	508
	Collateralized Debt Obligations (CDO's)	509
17.12	Summary	511
	Exercises	512
	References	513

**Part IV: GENERAL CASH FLOW STREAMS**

<b>Chapter 18 OPTIMAL PORTFOLIO GROWTH</b>		517
18.1	The Investment Wheel	517
	Analysis of the Wheel	519
18.2	The Log Utility Approach to Growth	519
	Log Utility Form	521
	Examples	521
18.3	Properties of the Log-Optimal Strategy*	525
18.4	Alternative Approaches*	526
	Other Utility	526
18.5	Continuous-Time Growth	528
	Dynamics of Several Stocks	528
	Portfolio Dynamics	529
	Implications for Growth	530
	The Portfolio of Maximum Growth Rate	530
18.6	The Feasible Region	531

The Efficient Frontier	531
Inclusion of a Risk-Free Asset	532
18.7 The Log-Optimal Pricing Formula*	536
Market Data	539
18.8 Log-Optimal Pricing and the Black–Scholes Equation*	540
18.9 Summary	541
Exercises	542
References	546
<b>Chapter 19 GENERAL INVESTMENT EVALUATION</b>	547
19.1 General Present Value	547
Projects and Opportunities	548
19.2 Multiperiod Securities*	548
Assets	549
Portfolio Strategies	549
Arbitrage	550
Short-Term Risk-Free Rates	550
19.3 Risk-Neutral Pricing	550
19.4 Optimal Pricing	552
The Single-Period Problem	552
Applications	553
19.5 The Double Lattice	555
19.6 Pricing in a Double Lattice	557
19.7 Investments with Private Uncertainty	560
General Approach	562
19.8 Buying Price Analysis	566
Certainty Equivalent and Exponential Utility	567
Sequential Calculation of CE	568
Multiperiod Case	569
General Approach	570
19.9 Pricing Axioms for Continuous Time	572
Option Formula	575
Risk-Neutral Form	575
Alternative Forms	575
19.10 Summary	576
Exercises	576
References	578
<b>Appendix A BASIC PROBABILITY THEORY</b>	579
A.1 General Concepts	579
A.2 Normal Random Variables	580
A.3 Lognormal Random Variables	581

<b>Appendix B CALCULUS AND OPTIMIZATION</b>	<b>583</b>
B.1 Functions	583
B.2 Differential Calculus	584
B.3 Optimization	585
<b>ANSWERS TO EXERCISES</b>	<b>588</b>
<b>INDEX</b>	<b>594</b>

# PREFACE

Investment is a fundamental component of modern life, reflected in all manner of economic activity. In practice, investment is generally carried out by processes facilitated by banks, mutual funds, brokers, and markets and governed by rules and protocols. These practicalities, together with the underlying investment motivation, comprise the related subjects of finance and investment. This overall field has recently expanded enormously, in terms of sheer volume but also in terms of the underlying theoretical structure. Recent developments in investment theory are being infused into university classrooms, into financial service organizations, into business ventures, and into the awareness of many individual investors. This book is intended to be one instrument in that dissemination process.

The book endeavors to emphasize fundamental principles and to illustrate how these principles can be mastered and transformed into sound and practical solutions of actual investment problems. The book's organizational structure reflects this approach: the material covered in the chapters progresses from the simplest in concept to the more advanced. Particular financial products and investment problems are treated, for the most part, in the order that they fall along this line of conceptual progression, their analyses serving to illustrate concepts as well as to describe particular features of the investment environment.

The book is designed for individuals who have a technical background roughly equivalent to a bachelor's degree in engineering, mathematics, or science; or who have some familiarity with basic mathematics. The language of investment science is largely mathematical, and some aspects of the subject can be expressed only in mathematical terms. The mathematics used in this book, however, is not complex—for example, only elementary portions of calculus are required—but the reader should be comfortable with the use of mathematics as a method of deduction and problem solving. Such readers will be able to leverage their technical backgrounds to accelerate and deepen their study.

Actually, the book can be read at several levels, requiring different degrees of mathematical sophistication and having different scopes of study. A simple road map to these different levels is coded into the typography of the text. Some section and subsection titles are set with an ending star as, for example, "2.6 Applications and Extensions.\*" The star indicates that the section or subsection is special: the material may be somewhat tangential or of higher mathematical level than elsewhere and can be skipped at first reading. This coding scheme is only approximate; the text itself often explains what is ahead in each section and gives guidelines on how the reader may wish to proceed.

The end-of-chapter exercises are an important part of the text, and readers should attempt several exercises in each chapter. The exercises are also coded: an exercise marked  $\diamond$  is mathematically more difficult than the average exercise; an exercise marked  $\oplus$  requires numerical computation (usually with a spreadsheet program).

Since publication of the first edition of this textbook, the subject of investment as a practical field and as an academic specialty has been extremely vibrant and innovative, with great interplay between theory and application, each motivating the other. As appropriate, much of this work is based on the fundamental concepts of CAPM and derivative theory, expanded and modified to address issues of portfolio design and risk management.

The real world of finance greatly tested the scope of traditional foundations. Issues of risk management, especially, became overwhelmingly important, as witnessed by the failure of some large banks and the high volatility and heavy losses in the stock market. Existing theory and its application methodologies, although sound, were not comprehensive enough. An early approach developed to measure the “riskiness” of institutions such as banks was *value at risk*, which by a single number quantifies a portfolio’s risk of loss. This measure has been widely accepted and indeed explicitly used for formal regulation of banks. The idea was studied by the academic community where variations such as *conditional value at risk* were proposed which have some theoretical and practical advantages over value at risk. Risk measures comprise the subject of Chapter 10.

Beyond simply measuring risk, credit derivatives were established that, like an insurance policy, protect the holder of a risky bond against default by the issuing entity. New theory was developed to price these credit derivatives. Much remains to be done in this area with regard to new products and theory. Credit risk is the subject of Chapter 17.

Another topic new to this edition is *projection pricing*, which relates to the common practice of applying the CAPM to price an asset that is not yet in the market. It is shown that this price can be found in other ways, one of which is related to the common everyday practice of pricing an asset by comparison with similar assets. Other new topics include a more comprehensive study of the effect of parameter estimation errors and how to minimize their negative impact, the “volatility smile” associated with option prices, and a simple axiomatic approach to valuation that unifies much of derivative pricing. The final chapter includes an extension of both CAPM and the Black–Scholes equation that prices continuous-time assets that are not derivatives. In addition, dozens of new end-of-chapter exercises (with answers to half of them) are included. Throughout, emphasis remains on a combination of clarity, intuition, and a modest level of mathematical rigor. Indeed some material has been modified to include new intuition to important theory.

The preparation of this second edition was a major project, one in which many people helped by their reading and critiquing of chapters. In this regard I wish to thank Giles Auchmuty, B. Ross Barmish, Xuedong He, Robert Kohn, Siu-Tang Leung, James Ligon and Frank Morgan. For special help with specific technical issues I wish to thank Samuel Chiu, Darrell Duffie, Daniel Gabay, Kay Giesecke, Marius Holtan, Daniel Kuhn, Robert Luenberger, Paul McEntire, James Primbs, Stan Uryasev, and Peter Woehmann. I am extremely grateful for students and graduates who helped in

this endeavor by reading chapters and helping devise new exercises. These include, Jose Blanchet, Naveed Chehraz, Ioannis Giannakakis, Supakorn Mudchanatongsuk, Ali Nouri, Dan Osborn, Wilfred Wong, and Li Xu. Finally, I (once again) thank my wife, Nancy, for her support during the long hours of manuscript development.

*DAVID G. LUENBERGER*

March 2013



# INTRODUCTION

**T**raditionally, investment is defined as the current commitment of resources in order to achieve later benefits. If resources and benefits take the form of money, investment is the present commitment of money for the purpose of receiving (hopefully more) money later. In some cases, such as the purchase of a bank certificate of deposit, the amount of money to be obtained later is known exactly. However, in most situations the amount of money to be obtained later is uncertain.

There is also a broader viewpoint of investment—based on flows of expenditures and receipts spanning a period of time. From this viewpoint, the objective of investment is to tailor the pattern of these flows over time to be as desirable as possible. When expenditures and receipts are denominated in cash, the net receipts at any time period are termed **cash flow**, and the series of flows over several periods is termed a **cash flow stream**. The investment objective is that of tailoring this cash flow stream to be more desirable than it would be otherwise. For example, by taking out a loan, it may be possible to exchange a large negative cash flow next month for a series of smaller negative cash flows over several months, and this alternative cash flow stream may be preferred to the original one. Often future cash flows have a degree of uncertainty, and part of the design, or tailoring, of a cash flow stream may be concerned with controlling that uncertainty, perhaps reducing the level of risk. This broader definition of investment, as tailoring a pattern of cash flows, encompasses the wide assortment of financial activities more fully than the traditional view. It is this broader interpretation that guides the presentation of this book.

**Investment science** is the application of scientific tools to investments. The scientific tools used are primarily mathematical, but only a modest level of mathematics

is required to understand the primary concepts discussed in this book. The purpose of this book is to convey both the principles of investment science and an understanding of how these principles can be used in practice to make calculations that lead to good investment decisions.

There is also an art to investment. Part of this art is knowing what to analyze and how to go about it. This part of the art can be enhanced by studying the material in this book. However, there is also an intuitive art of being able to evaluate an investment from an assortment of qualitative information, such as the personality characteristics of the people involved (the principals), whether a proposed new product will sell well, and so forth. This part of the art is not treated explicitly in this book, although the reader will gain some appreciation for just what this art entails.

## 1.1 Cash Flows

According to the broad interpretation, an investment is defined in terms of its resulting cash flow sequence—the amounts of money that will flow to and from an investor over time. Usually these cash flows (either positive or negative) occur at known specific dates, such as at the end of each quarter of a year or at the end of each year. The stream can then be described by listing the flow at each of these times. This is simplest if the flows are known deterministically, as in bank interest receipts or mortgage payments. In such cases the stream can be described by a series of numbers. For example, if the basic time period is taken as one year, one possible stream over a year, from beginning to end, is  $(-1, 1.2)$ , corresponding to an initial payment (the investment) of \$1 at the beginning of the year and the receipt of \$1.20 a year later. An investment over four years might be  $(-1, .10, .10, .10, 1.10)$ , where an initial investment of \$1 leads to payment of \$.10 at the end of each year for three years and then a final payment of \$1.10. Note that for a span of one year, two cash flow numbers are specified—one at the beginning and one at the end. Likewise, the four-year example involves five cash flow numbers.

Cash flow streams can also be represented in diagram form, as illustrated in Figure 1.1. In such a figure a time axis is drawn and a cash flow at a particular time is indicated by a vertical line at that time, the length of the line being proportional to the magnitude of the flow.



**FIGURE 1.1 Cash flow stream.** The cash flow stream of an investment can be represented by a diagram. In the example shown, the cash flows occur periodically. The first of these flows is negative, representing a cash outlay, and the subsequent flows are all positive.

If the magnitudes of some future cash flows are uncertain (as is frequently the case), a more complex representation of a cash flow stream must be employed. There are a few different techniques for doing this, and they are presented later in the book. But whether or not uncertainty is present, investments are described in terms of cash flow streams.

A diversity of investment issues can be stated in terms of cash flow streams, such as the following: Which of two cash flow streams is most preferable? How much would I be willing to pay to own a given stream? Are two streams together worth more to me than the sum of their individual values? If I can purchase a share of a stream, how much should I purchase? Given a collection of available cash flow streams, what is the most favorable combination of them?

Other more complex questions also arise. For example, sometimes the timing of all cash flows is not fixed, but can be influenced by the investor. If I purchase stock in a company, I have a negative cash flow initially, corresponding to my purchase payment; while I hold the stock, I perhaps receive dividends (relatively small positive cash flows) on a regular basis; finally, when I sell the stock, I obtain a major cash flow. However, the time of the last cash flow is not predetermined; I am free to choose it. Indeed, investments sometimes can be actively managed to influence both the amounts and the timing of all cash flows. For example, if I purchase a gold mine as an investment, I can decide how to mine it and thereby influence the cash flow every year. Determination of suitable management strategies is also part of investment science.

The view of investment science as the tailoring of cash flow streams gives the subject wide application. For individuals it applies to personal investment decisions, such as deciding on a home mortgage or planning for retirement. It also applies to business decisions, such as whether to invest in product development, whether to build a new manufacturing plant, and how to manage cash resources. Finally, it applies to government decisions, such as whether to build a dam or change the tax rate. Investment science guides us in the process of combining stocks, bonds, and other investment products into an overall package that has desirable properties. This process enhances total productivity by converting projects that in isolation may be too risky into members of attractive combinations.

## 1.2 Investments and Markets

At its root, investment analysis is a process of examining alternatives and deciding which alternative is most preferable. In this respect investment analysis is similar to the analysis of other decisions—operating a production facility, designing a building, planning a trip, or formulating an advertising campaign. Indeed, much of investment science relies on the same general tools used for analysis of these other decisions.

Investment problems differ from other decision problems in an important respect, however. Most investments are carried out within the framework of a financial market, and these markets provide alternatives not found in other decision situations. This structure is what makes investment analysis unique and unusually powerful.

## The Comparison Principle

Financial markets simplify decision making through a concept that we term the **comparison principle**. To introduce this principle, consider the following hypothetical situation.

Your uncle offers you a special investment. If you give him \$100 now, he will repay you \$110 in one year. His repayment is fully guaranteed by a trust fund of U.S. Treasury securities, and hence there is virtually no risk to the investment. Also, there is no moral or personal obligation to make this investment. You can either accept the offer or not. What should you do?

To analyze this situation, you would certainly note that the investment offers 10% interest, and you could compare this rate with the prevailing rate of interest that can be obtained elsewhere, say, at your local bank or from the U.S. Government through, for example, a Treasury bill. If the prevailing interest rate were only 7%, you would probably invest in this special offer by your uncle (assuming you have the cash to invest). If on the other hand the prevailing interest rate were 12%, you would surely decline the offer. From a pure investment viewpoint you can evaluate this opportunity very easily without engaging in deep reflection or mathematical analysis. If the investment offers a rate above normal, you accept; if it offers a rate below normal, you decline.

This analysis is an example of the comparison principle. You evaluate the investment by comparing it with other investments available in the financial market. The financial market provides a basis for comparison.

If, on the other hand, your uncle offers to sell you a family portrait whose value is largely sentimental, an outside comparison is not available. You must decide whether, to you, the portrait is worth his asking price.

## Arbitrage

When two similar investment alternatives are both available in the market, a principle stronger than the comparison principle holds. For example, consider (idealized) banks that offer to loan money or accept deposits at the same rate of interest. Suppose that the rate used at one bank for loans and deposits is 10% and at another bank the rate is 12%. You could go to the first bank and borrow, say, \$10,000 at 10% and then deposit that \$10,000 in the second bank at 12%. In one year you would earn 2% of \$10,000, which is \$200, without investing any cash of your own. This is a form of **arbitrage**—earning money without investing anything. Presumably, you could even make more money by running your scheme at a higher level. It should be clear that this kind of thing does not occur—at least not very often. The interest rates in the two banks would soon equalize.

The example of the two banks assumed that the interest rate for loans and the interest rate paid for deposits were equal within any one bank. Generally, of course, there is a difference in these rates. However, in markets of high volume, such as the

markets for U.S. Treasury securities, the difference between the buying price and the selling price is small. Therefore two different securities with identical properties must have approximately the same price—otherwise there would be an arbitrage opportunity.

Often it is assumed, for purposes of analysis, that **no arbitrage** opportunity exists. This is the no-arbitrage assumption.

Ruling out the possibility of arbitrage is a simple idea, but it has profound consequences. We shall find that the principle of no arbitrage implies that pricing relations are linear, that stock prices must satisfy certain relations, and that the prices of derivative securities, such as options and futures, can be determined analytically. This one principle, based on the existence of well-developed markets, permeates a good portion of modern investment science.

## Dynamics

Another important feature of financial markets is that they are dynamic, in the sense that the same or similar financial instruments are traded on a continuing basis. This means that the future price of an asset is not regarded as a single number, but rather as a process moving in time and subject to uncertainty. An important part of the analysis of an investment situation is the characterization of this process.

There are a few standard frameworks that are used to represent price processes. These include binomial lattice models, difference equation models, and differential equation models, all of which are discussed in this text. Typically, a record of past prices and other information are used to specify the parameters of such a model.

Because markets are dynamic, investment is itself dynamic—the value of an investment changes with time, and the quality of portfolios may change. Once this dynamic character is understood and formalized, it is possible to structure investments to take advantage of their dynamic nature so that the overall portfolio quality is maintained.

## Risk Aversion

Another principle of investment science is risk aversion. Suppose two possible investments have the same cost, and both are expected to return the same amount (somewhat greater than the initial cost), where the term *expected* is defined in a probabilistic sense (explained in Chapter 6). However, the return is certain for one of these investments, and the return on the other, while statistically independent of other assets, is uncertain. Individuals seeking investment rather than outright speculation will elect the first (certain) alternative over the second (risky) alternative. This is the risk aversion principle.

This risk aversion principle can be formalized (and made analytical) in a few different ways, which are discussed in later chapters. Once a formalism is established, the risk aversion principle can be used to help analyze many investment alternatives.

One way that the risk aversion principle is formalized is through **mean-variance analysis**. In this approach, the uncertainty of the return on a complete portfolio is characterized by just two quantities: the mean value of the return and the variance of the return. The risk aversion principle then says that if several possible portfolios have the same mean but different variances, a rational (risk-averse) investor will select the one that has the smallest variance.

This mean–variance method of formalizing risk is the basis for the most well-known method of quantitative portfolio analysis, which was pioneered by Harry Markowitz (who won the Nobel prize in economics for his work). This approach leads to a comprehensive theory of investment and is widely considered to be the foundation for modern portfolio theory. We discuss this important theory in Chapter 6.

A more general way to formalize the risk aversion principle is through the introduction of individual **utility functions**. This approach is presented in Chapter 11.

Later, in Chapter 18, we find that risk aversion takes on a new character when investments are made repeatedly over time. In fact, short-term variance will be found to be *good*, not bad. This is one of the surprising conclusions of the comprehensive view of investment represented by investment science.

## 1.3 Typical Investment Problems

Every investment problem has unique features, but many fit into a few broad categories or types. We briefly outline some of the most important problem types here. Fuller descriptions of these general types and more specific examples appear in the relevant chapters.

### Pricing

Let us go back to our very first example of an investment situation, the first offer from your uncle, but now let us turn it around. Imagine that there is an investment opportunity that will pay exactly \$110 at the end of one year. We ask: How much is this investment worth today? In other words, what is the appropriate price of this investment, given the overall financial environment?

If the current interest rate for one-year investments is 10%, then this investment should have a price of exactly \$100. In that case, the \$110 paid at the end of the year would correspond to a rate of return of 10%. If the current interest rate for one-year investments is less than 10%, then the price of this investment would be somewhat greater than \$100. In general, if the interest rate is  $r$  (expressed as a decimal, such as  $r = .10$ ), then the price of an investment that pays  $X$  after one year should be  $X/(1+r)$ .

We determined the price by a simple application of the comparison principle. This investment can be directly compared with one of investing money in a one-year certificate of deposit (or one-year Treasury bill), and hence it must bear the same effective interest rate.

This interest rate example is a simple example of the general pricing problem: Given an investment with known payoff characteristics (which may be random), what is the reasonable price; or, equivalently, what price is consistent with the other securities that are available? We shall encounter this problem in many contexts. For example, early in our study we shall determine the appropriate price of a bond. Later we shall compute the appropriate price of a share of stock with random return characteristics. Still later we shall compute suitable prices of more complicated securities, such as futures and options. Indeed, the pricing problem is one of the basic problems of modern investment science and has obvious practical applications.

As in the simple interest rate example, the pricing problem is usually solved by use of the comparison principle. In most instances, however, the application of that principle is not as simple and obvious as in this example. Clever arguments have been devised to show how a complex investment can be separated into parts, each of which can be compared with other investments whose prices are known. Nevertheless, whether by a simple or a complex argument, comparison is the basis for the solution of many pricing problems.

## Hedging

**Hedging** is the process of reducing the financial risks that either arise in the course of normal business operations or are associated with investments. Hedging is one of the most important uses of financial markets, and is an essential part of modern industrial activity. One form of hedging is **insurance** where, by paying a fixed amount (a **premium**), you can protect yourself against certain specified possible losses—such as losses due to fire, theft, or even adverse price changes—by arranging to be paid compensation for the losses you incur. More general hedging can arise in the following way. Imagine a large bakery. This bakery will purchase flour (made from wheat and other ingredients) and transform these ingredients into baked goods, such as bread. Suppose the bakery wins a contract to supply a large quantity of bread to another company over the next year at a fixed price. The bakery is happy to win the contract, but now faces risk with respect to flour prices. The bakery will not immediately purchase all the flour needed to satisfy the contract, but will instead purchase flour as needed during the year. Therefore, if the price of flour should increase part way during the year, the bakery will be forced to pay more to satisfy the needs of the contract and, hence, will have a lower profit. In a sense the bakery is at the mercy of the flour market. If the flour price goes up, the bakery will make less profit, perhaps even losing money on the contract. If the flour price goes down, the bakery will make even more money than anticipated.

The bakery wishes to be in the baking business, not in the flour speculation business. It wants to eliminate the risk associated with flour costs and concentrate on baking. It can do this by obtaining an appropriate number of wheat futures contracts in the futures market. Such a contract has small or zero initial cash outlay and at a set future date gives a profit (or loss) equal to the amount that wheat prices have changed since entering the contract. The price of flour is closely tied to the price of wheat, so if the price of flour should go up, the value of a wheat futures contract will go up by

a somewhat comparable amount. Hence the net effect to the bakery—the profit from the wheat futures contracts together with the change in the cost of flour—is nearly zero.

There are many other examples of business risks that can be reduced by hedging. And there are many ways that hedging can be carried out: through futures contracts, options, and other special arrangements. Indeed, the major use, by far, of these financial instruments is for hedging—not for speculation.

## Risk Assessment and Management

Frequently, the issue of risk dominates the appraisal of an investment or a portfolio. This occurs when an individual makes a very large investment and there is a possibility (seemingly remote) of a catastrophic loss. For example, when purchasing a house, a young couple may focus on unique physical risks (such as whether the house is on solid soil) or on financial risks (such as whether the housing market will fall). Coming to terms with these risks and possibly insuring against them may be more important than further negotiation involving, say, 10% of the sales price. Most large financial institutions, and the government regulators that oversee them, rely on quantitative methods for assessing potentially devastating risks. This general field, referred to as *risk management*, is an important part of investment science. A complication of this subject is that risk assessments are based on probabilistic models of possibilities, and it is difficult to determine the reliability of the underlying models.

A special kind of risk, *counterparty risk*, arises when entering a contract. This is the risk that the other party may default. For example, individuals face this risk whenever they deposit money in a bank, for the bank may fail. And financial institutions face this sort of risk when dealing with other institutions. Fortunately, insurance is often available either from a government agency or in the form of special securities designed for this purpose. As an example, a bank that owns the bonds of a municipality may purchase a special **credit default swap** that pays a predetermined amount to the bank if the counterparty municipality defaults on the bond. Insurance products of this sort are themselves derivative securities, and their prices are analyzed by the general theories and methods of derivatives.

## Pure Investment

Pure investment refers to the objective of obtaining increased future return for present allocation of capital. This is the motivation underlying most individual investments in the stock market, for example. The investment problem arising from this motivation is referred to as the **portfolio selection problem**, since the real issue is to determine where to invest available capital.

Most approaches to the pure investment problem rely on the risk aversion principle, for in this problem one must carefully assess one's preferences, deciding how to balance risk and expected reward. There is not a unique solution. Judgment

and taste are important, which is evidenced by the vast amount of literature and advice directed each year to helping individuals find solutions to this problem.

The pure investment problem also characterizes the activities of a profit-seeking firm which, after all, takes existing capital and transforms it, through investment—in equipment, people, and operations—into profit. Hence the methods developed for analyzing pure investment problems can be used to analyze potential projects within firms, the overall financial structure of a firm, and even mergers of firms.

## Other Problems

Investment problems do not always take the special shapes outlined in the preceding categories. A hedging problem may contain an element of pure investment, and conversely an investment may be tempered with a degree of hedging. Fortunately, the same principles of analysis are applicable to such combinations.

One type of problem that occurs frequently is a combined consumption-investment problem. For example, a married couple at retirement, living off the income from their investments, will most likely invest differently than a young couple investing for growth of capital. The requirement for income changes the nature of the investment problem. Likewise, the management of an endowment for a public enterprise, such as a university must consider growth objectives as well as consumption-like objectives associated with the current operations of the enterprise.

We shall also find that the framework of an investment problem is shaped by the formal methods used to treat it. Once we have logical methods for representing investment issues, new problems suggest themselves. As we progress through the book we shall uncover additional problems and obtain a deeper appreciation for the simple outlines given here.

## 1.4 Organization of the Book

The organization of this book reflects the notion that investment science is the study of how to tailor cash flow streams. Indeed, the cash flow viewpoint leads to a natural partition of the subject into four main parts, as follows.

### Deterministic Cash Flow Streams

The simplest cash flow streams are those that are deterministic (that is, not random, but definite). The first part of the book treats these. Such cash flows can be represented by sequences such as  $(-1, 0, 3)$ , as discussed earlier. Investments of this type, either with one or with several periods, are analyzed mainly with various concepts of interest rate. Accordingly, interest rate theory is emphasized in this first part of the book. This theory provides a basis for a fairly deep understanding of investment and a framework for addressing a wide variety of important and interesting problems.

## Single-Period Random Cash Flow Streams

The second level of complexity in cash flow streams is associated with streams having only a single period, with beginning and ending flows, but with the magnitude of the second flow being uncertain. Such a situation occurs when a stock is purchased at the beginning of the year and sold at the end of the year. The amount received at the end of the year is not known in advance and, hence, must be considered uncertain. This level of complexity captures the essence of many investment situations.

In order to analyze cash flows of this kind, one must have a formal description of **uncertain returns**. There are several such descriptions (all based on probability theory), and we shall study the main ones, the simplest being the **mean–variance** description. One must also have a formal description of how individuals assess uncertain returns. We shall consider such assessment methods, starting with mean–variance analysis. These single-period uncertain cash flow situations are the subject of the second part of the book.

## Derivative Assets

The third level of complexity in cash flow streams involves streams that have random flows at each of several time points, but where the asset producing a stream is functionally related to another asset whose price characteristics are known.

An asset whose cash flow values depend functionally on another asset is termed a **derivative asset**. A good example is a stock option. To describe such an option, suppose that I own 100 shares of stock in company A. This asset, the 100 shares, is a **basic asset**. Now suppose that I have granted you the right (but not the obligation) to buy, at say \$54 per share, all 100 of my shares in three months. This right is a call option on 100 shares of stock in company A. This option is an asset; it has value, and that value may change with time. It is, however, a derivative of the stock of company A because the value of the option depends on the price of the stock. If the stock price goes up, the option value also goes up. Other derivative assets include futures contracts, other kinds of options, and various other financial contracts. One example seen by many home buyers is the adjustable-rate mortgage, which periodically adjusts interest payments according to an interest rate index. Such a mortgage is a derivative of the securities that determine the interest rate index.

The third part of the book is devoted to these derivative assets. Analysis of these assets is often simpler than that for assets with general multiperiod uncertain cash flows because properties of a derivative can be traced back to the underlying basic asset. The study of derivative assets, however, is an important and lively aspect of investment science, one for which strong theoretical results can be derived and important numerical quantities, such as implied prices, can be obtained.

## General Cash Flow Streams

Finally, the fourth part of the book is devoted to cash flow streams with uncertain cash flows at many different times—flows that are not functionally related to other assets. As can be expected, this final level of complexity is the most difficult part of the subject, but also the one that is the most important. The cash flow streams encountered in most investments have this general form.

The methods of this part of the book build on those of earlier parts, but new concepts are added. The fact that the mix of investments—the portfolio structure—can be changed as time progresses, depending on what has happened to that point, leads to new phenomena and new opportunities. For example, the growth rate of a portfolio can be enhanced by employing suitable reinvestment strategies. This part of the book represents some of the newest aspects of the field.

The theory of general cash flows—random with many periods—enables us to address many practical investment problems arising in business situations. For example, a firm may consider a project with associated cash flows that are random but possibly influenced by management actions. These actions often can be characterized as options (such as an option to modify a production process or to vary pricing). Since the project value in this case is the value of a “real” asset rather than of a purely “paper” financial asset, these situations are frequently referred to as **real options**. In some cases the standard theory of derivatives can be used to find this value. However, application of the standard theory requires that the cash flows be derivative of some marketed underlying security (which at each moment can provide a perfect hedge). In many cases, no such underlying security exists, in which case other approaches must be devised. The material of the early chapters of the text can be appropriately adapted to these situations.

Investment science is a practical science; and because its main core is built on a few simple principles, it can be easily learned and fruitfully applied to interesting and important problems. It is also an evolving science, which is expanding rapidly. Perhaps the reader, armed with a basic understanding of the field, will contribute to this evolution through either theory or application.



## PART I

# DETERMINISTIC CASH FLOW STREAMS





# 2

## THE BASIC THEORY OF INTEREST

Interest is frequently called the time value of money, and the next few chapters explore the structure and implications of this value. In this first chapter on the subject, we outline the basic elements of interest rate theory, showing that the theory can be translated into analytic form and thus used as a basis for making intelligent investment decisions.

### 2.1 Principal and Interest

The basic idea of interest is quite familiar. If you invest \$1.00 in a bank account that pays 8% interest per year, then at the end of 1 year you will have in your account the **principal** (your original amount) of \$1.00 plus **interest** of \$.08 for a total of \$1.08. If you invest a larger amount, say  $A$  dollars, then at the end of the year your account will have grown to  $A \times 1.08$  dollars. In general, if the interest rate is  $r$ , expressed as a decimal, then your initial investment would be multiplied by  $(1 + r)$  after 1 year.

#### Simple Interest

Under a **simple interest** rule, money invested for a period different from 1 year accumulates interest proportional to the total time of the investment. So after 2 years, the total interest due is  $2r$  times the original investment, and so forth. In other words, the investment produces interest equal to  $r$  times the original investment every year.

Usually partial years are treated in a proportional manner; that is, after a fraction  $f$  of 1 year, interest of  $rf$  times the original investment is earned.

The general rule for simple interest is that if an amount  $A$  is left in an account at simple interest  $r$ , the total value after  $n$  years is

$$V = (1 + rn)A.$$

If the proportional rule holds for fractional years, then after any time  $t$  (measured in years), the account value is

$$V = (1 + rt)A.$$

The account grows **linearly** with time. As shown in the preceding formula, the account value at any time is just the sum of the original amount (the principal) and the accumulated interest, which is proportional to time.

## Compound Interest

Most bank accounts and loans employ some form of compounding—producing compound interest. Again, consider an account that pays interest at a rate of  $r$  per year. If interest is compounded yearly, then after 1 year, the first year's interest is added to the original principal to define a larger principal base for the second year. Thus during the second year, the account earns *interest on interest*. This is the compounding effect, which is continued year after year.

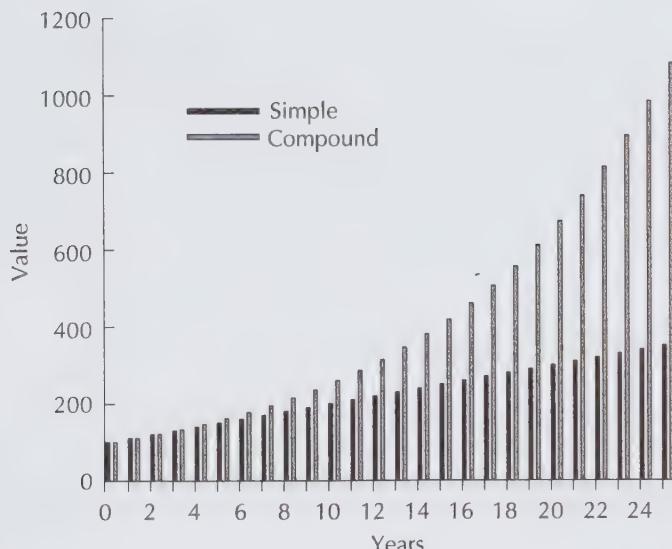
Under yearly compounding, money left in an account is multiplied by  $(1 + r)$  after 1 year. After the second year, it grows by another factor of  $(1 + r)$  to  $(1 + r)^2$ . After  $n$  years, such an account will grow to  $(1 + r)^n$  times its original value, and this is the analytic expression for the account growth under **compound interest**. This expression is said to exhibit **geometric growth** because of its  $n$ th-power form.

As  $n$  increases, the growth due to compounding can be substantial. For example, Figure 2.1 shows a graph of a \$100 investment over time when it earns 10% interest under simple and compound interest rules. The figure shows the characteristic shapes of linear growth for simple interest and of accelerated upward growth for compound interest. Note that under compounding, the value doubles in about 7 years.

There is a cute little rule that can be used to estimate the effect of interest compounding.

**The seven-ten rule** *Money invested at 7% per year doubles in approximately 10 years. Also, money invested at 10% per year doubles in approximately 7 years.*

(More exactly, at 7% and 10 years, an account increases by a factor of 1.97, whereas at 10% and 7 years it increases by a factor of 1.95.) The rule can be generalized, and slightly improved, to state that, for interest rates less than about 20%, the doubling time is approximately  $72/i$ , where  $i$  is the interest rate expressed as a percentage (that is, 10% interest corresponds to  $i = 10$ ). (See Exercise 2.)



**FIGURE 2.1 Simple and compound interest.** Simple interest leads to linear growth over time, whereas compound interest leads to an accelerated increase defined by geometric growth. The figure shows both cases for an interest rate of 10%.

## Compounding at Various Intervals

In the preceding discussion, interest was calculated at the end of each year and paid to the account at that time. Most banks now calculate and pay interest more frequently—quarterly, monthly, or in some cases daily. This more frequent compounding raises the effective yearly rate. In this situation, it is traditional to still quote the interest rate on a yearly basis, but then apply the appropriate proportion of that interest rate over each compounding period. For example, consider quarterly compounding. Quarterly compounding at an interest rate of  $r$  per year means that an interest rate of  $r/4$  is applied every quarter. Hence money left in the bank for 1 quarter will grow by a factor of  $1 + (r/4)$  during that quarter. If the money is left in for another quarter, then that new amount will grow by another factor of  $1 + (r/4)$ . After 1 year the account will have grown by the compound factor of  $[1 + (r/4)]^4$ . For any  $r > 0$ , it holds that  $[1 + (r/4)]^4 > 1 + r$ . Hence at the same annual rate, the amount in the bank account after 4 quarters of compounding is greater than the amount after 1 year without compounding.

The effect of compounding on yearly growth is highlighted by stating an **effective interest rate**, which is the equivalent yearly interest rate that would produce the same result after 1 year without compounding. For example, an annual rate of 8% compounded quarterly will produce an increase of  $(1.02)^4 = 1.0824$ ; hence the effective interest rate is 8.24%. The basic yearly rate (8% in this example) is termed the **nominal rate**.

Compounding can be carried out with any frequency. The general method is that a year is divided into a fixed number of equally spaced periods—say  $m$  periods. (In the case of monthly compounding the periods are not quite equal, but we shall ignore that here and regard monthly compounding as simply setting  $m = 12$ .) The interest rate for each of the  $m$  periods is thus  $r/m$ , where  $r$  is the nominal annual rate. The account grows by  $1 + (r/m)$  during 1 period. After  $k$  periods, the growth is  $[1 + (r/m)]^k$ , and hence after a full year of  $m$  periods it is  $[1 + (r/m)]^m$ . The effective interest rate is the number  $r'$  that satisfies  $1 + r' = [1 + (r/m)]^m$ .

## Continuous Compounding

We can imagine dividing the year into smaller and smaller periods, and thereby apply compounding monthly, weekly, daily, or even every minute or second. This leads to the idea of continuous compounding. We can determine the effect of continuous compounding by considering the limit of ordinary compounding as the number  $m$  of periods in a year goes to infinity. To determine the yearly effect of this continuous compounding we use the fact that

$$\lim_{m \rightarrow \infty} [1 + (r/m)]^m = e^r$$

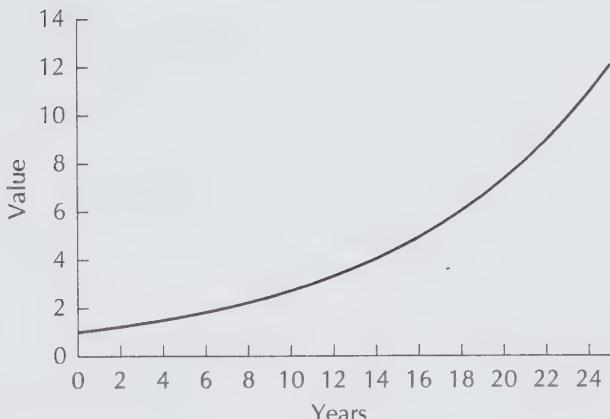
where  $e = 2.7818\dots$  is the base of the natural logarithm. The effective rate of interest  $r'$  is the value satisfying  $1 + r' = e^r$ . If the nominal interest rate is 8% per year, then with continuous compounding the growth would be  $e^{.08} = 1.0833$ , and hence the effective interest rate is 8.33%. (Recall that quarterly compounding produces an effective rate of 8.24%.) Table 2.1 shows the effect of continuous compounding for various nominal rates. Note that as the nominal rate increases, the compounding effect becomes more dramatic.

We can also calculate how much an account will have grown after any arbitrary length of time. We denote time by the variable  $t$ , measured in years. Thus  $t = 1$  corresponds to 1 year, and  $t = .25$  corresponds to 3 months. Select a time  $t$  and divide the year into a (large) number  $m$  of small periods, each of length  $1/m$ . Then  $t \approx k/m$  for some  $k$ , meaning that  $k$  periods approximately coincides with the time  $t$ . If  $m$  is very large, this approximation can be made very accurate. Therefore  $k \approx mt$ . Using

**TABLE 2.1**  
**CONTINUOUS COMPOUNDING**

	Interest rate (%)							
Nominal	1.00	5.00	10.00	20.00	30.00	50.00	75.00	100.00
Effective	1.01	5.13	10.52	22.14	34.99	64.87	111.70	171.83

The nominal interest rates in the top row correspond, under continuous compounding, to the effective rates shown in the second row. The increase due to compounding is quite dramatic at large nominal rates.



**FIGURE 2.2 Exponential growth curve; continuous compound growth.** Under continuous compounding at 10%, the value of \$1 doubles in about 7 years. In 20 years it grows by a factor of about 8.

the general formula for compounding, we know that the growth factor for  $k$  periods is

$$[1 + (r/m)]^k = [1 + (r/m)]^{mt} = \{[1 + (r/m)]^m\}^t \rightarrow e^{rt}$$

where that last expression is valid in the limit as  $m$  goes to infinity, corresponding to continuous compounding. Hence continuous compounding leads to the familiar **exponential growth** curve. Such a curve is shown in Figure 2.2 for a 10% nominal interest rate.

## Debt

We have examined how a single investment (say a bank deposit) grows over time due to interest compounding. It should be clear that exactly the same thing happens to debt. If I *borrow* money from the bank at an interest rate  $r$  and make no payments to the bank, then my debt increases according to the same formulas. Specifically, if my debt is compounded monthly, then after  $k$  months my debt will have grown by a factor of  $[1 + (r/12)]^k$ .

## Money Markets

Although we have treated interest as a given known value, in reality there are many different rates each day. Different rates apply to different circumstances, different user classes, and different periods. For example, U.S. Treasury bills are issued with maturities of 3 or 13, 26 or 52 weeks. U.S. Treasury notes are issued with maturities of 2, 3, 5, 7, and 10 years. U.S. Treasury bonds have maturity of 30 years. These all have their individual rates and are discussed more fully in Chapter 3. In addition, there

are bonds issued by corporations or municipalities, and these have various implied interest rates. There is also a whole host of rates applying to loans between banks, home mortgages, and business loans. Most rates are established by the forces of supply and demand in broad markets to which they apply. Many of these market rates are discussed more fully in Chapters 3 and 4. Not all interest rates are broad market rates. There may be private rates negotiated by two private parties. Or in the context of a firm, special rates may be established for internal transactions or for the purpose of evaluating projects, as discussed later in this chapter.

## 2.2 Present Value

The theme of the previous section is that money invested today leads to increased value in the future as a result of interest. The formulas of the previous section show how to determine this future value.

That whole set of concepts and formulas can be reversed in time to calculate the value that should be assigned now, in the present, to money that is to be received at a later time. This reversal is the essence of the extremely important concept of **present value**.

To introduce this concept, consider two situations: (1) you will receive \$110 in 1 year, (2) you receive \$100 now and deposit it in a bank account for 1 year at 10% interest. Clearly these situations are identical after 1 year—you will receive \$110. We can restate this equivalence by saying that \$110 received in 1 year is equivalent to the receipt of \$100 now when the interest rate is 10%. Or we say that the \$110 to be received in 1 year has a **present value** of \$100. In general, \$1 to be received a year in the future has a present value of  $\$1/(1+r)$ , where  $r$  is the interest rate.

A similar transformation applies to future obligations such as the repayment of debt. Suppose that, for some reason, you have an obligation to pay someone \$100 in exactly 1 year. This obligation can be regarded as a negative cash flow that occurs at the end of the year. To calculate the present value of this obligation, you determine how much money you would need *now* in order to cover the obligation. This is easy to determine. If the current yearly interest rate is  $r$ , you need  $\$100/(1+r)$ . If that amount of money is deposited in the bank now, it will grow to \$100 at the end of the year. You can then fully meet the obligation. The present value of the obligation is therefore  $100/(1+r)$ .

The process of evaluating future obligations as an equivalent present value is alternatively referred to as **discounting**. The present value of a future monetary amount is less than the face value of that amount, so the future value must be discounted to obtain the present value. The factor by which the future value must be discounted is called the **discount factor**. The 1-year discount factor is  $d_1 = 1/(1+r)$ , where  $r$  is the 1-year interest rate. So if an amount  $A$  is to be received in 1 year, the present value is the discounted amount  $d_1A$ .

The formula for present value depends on the interest rate that is available from a bank or other source. If that source quotes rates with compounding, then such a compound interest rate should be used in the calculation of present value. As an example, suppose that the annual interest rate  $r$  is compounded at the end of each

of  $m$  equal periods each year; and suppose that a cash payment of amount  $A$  will be received at the end of the  $k$ th period. Then the appropriate discount factor is

$$d_k = \frac{1}{[1 + (r/m)]^k}.$$

The present value of a payment of  $A$  to be received  $k$  periods in the future is  $d_k A$ .

## 2.3 Present and Future Values of Streams

The previous section studied the impact of interest on a single cash deposit or loan; that is, on a single cash flow. We now extend that discussion to the case where cash flows occur at several time periods, and hence constitute a cash flow stream or sequence. First we require a new concept.

### The Ideal Bank

When discussing cash flow streams, it is useful to define the notion of an **ideal bank**. An ideal bank applies the same rate of interest to both deposits and loans, and it has no service charges or transaction fees. Its interest rate applies equally to any size of principal, from 1 cent (or fraction thereof) to \$1 million (or even more). Furthermore, separate transactions in an account are completely additive in their effect on future balances.

Note that the definition of an ideal bank does *not* imply that interest rates for all transactions are identical. For example, a 2-year certificate of deposit (CD) might offer a higher rate than a 1-year CD. However, the 2-year CD must offer the same rate as a loan that is payable in 2 years.

If an ideal bank has an interest rate that is independent of the length of time for which it applies, and that interest is compounded according to normal rules, it is said to be a **constant ideal bank**. In the rest of this chapter, we always assume that interest rates are indeed constant.

The constant ideal bank is the reference point used to describe the outside financial market—the public market for money.

### Future Value

Now we return to the study of cash flow streams. Let us decide on a fixed time cycle for compounding (for example, yearly) and let a period be the length of this cycle. We assume that cash flows occur at the end of each period (although some flows might be zero). We shall take each cash flow and deposit it in a constant ideal bank as it arrives. (If the flow is negative, we cover it by taking out a loan.) Under the terms of a constant ideal bank, the final balance in our account can be found by combining the results of the individual flows. Explicitly, consider the cash flow stream  $(x_0, x_1, \dots, x_n)$ . At the end of  $n$  periods the initial cash flow  $x_0$  will have grown

to  $x_0(1+r)^n$ , where  $r$  is the interest rate *per period* (which is the yearly rate divided by the number of periods per year). The next cash flow,  $x_1$ , received after the first period, will at the final time have been in the account for only  $n-1$  periods, and hence it will have a value of  $x_1(1+r)^{n-1}$ . Likewise, the next flow  $x_2$  will collect interest during  $n-2$  periods and have value  $x_2(1+r)^{n-2}$ . The final flow  $x_n$  will not collect any interest, so will remain  $x_n$ . The total value at the end of  $n$  periods is therefore  $FV = x_0(1+r)^n + x_1(1+r)^{n-1} + \dots + x_n$ . To summarize:

**Future value of a stream** *Given a cash flow stream  $(x_0, x_1, \dots, x_n)$  and interest rate  $r$  each period, the future value of the stream is*

$$FV = x_0(1+r)^n + x_1(1+r)^{n-1} + \dots + x_n.$$

**Example 2.1 (A short stream)** Consider the cash flow stream  $(-2, 1, 1, 1)$  when the periods are years and the interest rate is 10%. The future value is

$$FV = -2 \times (1.1)^3 + 1 \times (1.1)^2 + 1 \times 1.1 + 1 = .648. \quad (2.1)$$

This formula for future value always uses the interest rate per period and assumes that interest rates are compounded each period.

## Present Value

The present value of a general cash flow stream—like the future value—can also be calculated by considering each flow element separately. Again consider the stream  $(x_0, x_1, \dots, x_n)$ . The present value of the first element  $x_0$  is just that value itself since no discounting is necessary. The present value of the flow  $x_1$  is  $x_1/(1+r)$ , because that flow must be discounted by one period. (Again the interest rate  $r$  is the per-period rate.) Continuing in this way, we find that the present value of the entire stream is  $PV = x_0 + x_1/(1+r) + x_2/(1+r)^2 + \dots + x_n/(1+r)^n$ . We summarize this important result as follows.

**Present value of a stream** *Given a cash flow stream  $(x_0, x_1, \dots, x_n)$  and an interest rate  $r$  per period, the present value of this cash flow stream is*

$$PV = x_0 + \frac{x_1}{1+r} + \frac{x_2}{(1+r)^2} + \dots + \frac{x_n}{(1+r)^n}. \quad (2.2)$$

**Example 2.2** Again consider the cash flow stream  $(-2, 1, 1, 1)$ . Using an interest rate of 10% we have

$$PV = -2 + \frac{1}{1.1} + \frac{1}{(1.1)^2} + \frac{1}{(1.1)^3} = .487.$$

The present value of a cash flow stream can be regarded as the present payment amount that is equivalent to the entire stream. Thus we can think of the entire stream as being replaced by a single flow at the initial time.

There is another way to interpret the formula for present value that is based on transforming the formula for future value. Future value is the amount of future payment that is equivalent to the entire stream. We can think of the stream as being transformed into that single cash flow at period  $n$ . The present value of this single equivalent flow is obtained by discounting it by  $(1 + r)^n$ . That is, the present value and the future value are related by

$$PV = \frac{FV}{(1 + r)^n}.$$

In the previous examples for the cash flow stream  $(-2, 1, 1, 1)$  we have  $.487 = PV = FV/(1.1)^3 = .648/1.331 = .487$ .

## Frequent and Continuous Compounding

Suppose that  $r$  is the nominal annual interest rate and interest is compounded at  $m$  equally spaced periods per year. Suppose that cash flows occur initially and at the end of each period for a total of  $n$  periods, forming a stream  $(x_0, x_1, \dots, x_n)$ . Then according to the preceding we have

$$PV = \sum_{k=0}^n \frac{x_k}{[1 + (r/m)]^k}.$$

Suppose now that the nominal interest rate  $r$  is compounded continuously and cash flows occur at times  $t_0, t_1, \dots, t_n$ . (We have  $t_k = k/m$  for the stream in the previous paragraph; but the more general situation is allowed here.) We denote the cash flow at time  $t_k$  by  $x(t_k)$ . In that case,

$$PV = \sum_{k=0}^n x(t_k) e^{-rt_k}.$$

This is the continuous compounding formula for present value.

## Present Value and an Ideal Bank

We know that an ideal bank can be used to change the pattern of a cash flow stream. For example, a 10% bank can change the stream  $(1, 0, 0)$  into the stream  $(0, 0, 1.21)$  by receiving a deposit of \$1 now and paying principal and interest of \$1.21 in 2 years. The bank can also work in a reverse fashion and transform the second stream into the first by issuing a loan for \$1 now.

In general, if an ideal bank can transform the stream  $(x_0, x_1, \dots, x_n)$  into the stream  $(y_0, y_1, \dots, y_n)$ , it can also transform in the reverse direction. Two streams that can be transformed into each other are said to be **equivalent streams**.

How can we tell whether two given streams are equivalent? The answer to this is the main theorem on present value.

**Main theorem on present value** *The cash flow streams  $\mathbf{x} = (x_0, x_1, \dots, x_n)$  and  $\mathbf{y} = (y_0, y_1, \dots, y_n)$  are equivalent for a constant ideal bank with interest rate  $r$  if and only if the present values of the two streams, evaluated at the bank's interest rate, are equal.*

**Proof:** Let  $v_{\mathbf{x}}$  and  $v_{\mathbf{y}}$  be the present values of the  $\mathbf{x}$  and  $\mathbf{y}$  streams, respectively. Then the  $\mathbf{x}$  stream is equivalent to the stream  $(v_{\mathbf{x}}, 0, 0, \dots, 0)$  and the  $\mathbf{y}$  stream is equivalent to the stream  $(v_{\mathbf{y}}, 0, 0, \dots, 0)$ .

It is clear that these two streams are equivalent if and only if  $v_{\mathbf{x}} = v_{\mathbf{y}}$ . Hence the original streams are equivalent if and only if  $v_{\mathbf{x}} = v_{\mathbf{y}}$ . ■

This result is important because it implies that present value is the only number needed to characterize a cash flow stream when an ideal bank is available. The stream can be transformed in a variety of ways by the bank, but the present value remains the same. So if someone offers you a cash flow stream, you only need to evaluate its corresponding present value, because you can then use the bank to tailor the stream with that present value to any shape you desire.

## 2.4 Internal Rate of Return

**Internal rate of return** is another important concept of cash flow analysis. It pertains specifically to the entire cash flow stream associated with an investment, not to a partial stream such as a cash flow at a single period. The streams to which this concept is applied typically have both negative and positive elements: the negative flows correspond to the payments that must be made; the positive flows to payments received. A simple example is the process of investing in a certificate of deposit for a fixed period of 1 year. Here there are two cash flow elements: the initial deposit or payment (a negative flow) and the final redemption (a positive flow).

Given a cash flow stream  $(x_0, x_1, \dots, x_n)$  associated with an investment, we write the present value formula

$$PV = \sum_{k=0}^n \frac{x_k}{(1+r)^k}.$$

If the investment that corresponds to this stream is constructed from a series of deposits and withdrawals from a constant ideal bank at interest rate  $r$ , then from the main theorem on present value of the previous section, PV would be zero. The idea behind internal rate of return is to turn the procedure around. Given a cash flow stream, we write the expression for present value and then find the value of  $r$  that renders this present value equal to zero. That value is called the internal rate of return because it is the interest rate implied by the internal structure of the cash flow stream. The idea can be applied to any series of cash flows.

The preliminary formal definition of the internal rate of return (IRR) is as follows:

**Internal rate of return** *Let  $(x_0, x_1, x_2, \dots, x_n)$  be a cash flow stream. Then the internal rate of return of this stream is a number  $r$  satisfying the equation*

$$0 = x_0 + \frac{x_1}{1+r} + \frac{x_2}{(1+r)^2} + \cdots + \frac{x_n}{(1+r)^n}. \quad (2.3)$$

*Equivalently, it is a number  $r$  satisfying  $1/(1+r) = c$  [that is,  $r = (1/c) - 1$ ], where  $c$  satisfies the polynomial equation*

$$0 = x_0 + x_1c + x_2c^2 + \cdots + x_nc^n. \quad (2.4)$$

We call this a preliminary definition because there may be ambiguity in the solution of the polynomial equation of degree  $n$ . We discuss this point shortly. First, however, let us illustrate the computation of the internal rate of return.

**Example 2.3 (The old stream)** Consider again the cash flow sequence  $(-2, 1, 1, 1)$  discussed earlier. The internal rate of return is found by solving the equation

$$0 = -2 + c + c^2 + c^3.$$

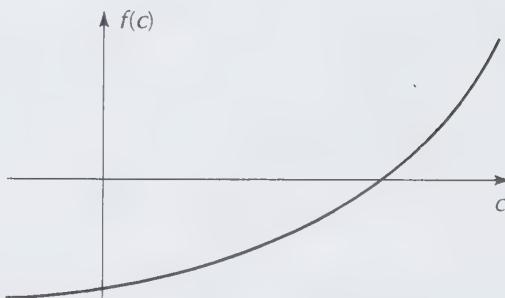
The positive solution can be found (by trial and error) to be  $c = .81$ , and thus  $\text{IRR} = (1/c) - 1 = .23$ .

Notice that the internal rate of return is defined without reference to a prevailing interest rate. It is determined entirely by the cash flows of the stream. This is the reason why it is called the *internal* rate of return; it is defined internally without reference to the external financial world. It is the rate that an ideal bank would have to apply to generate the given stream from an initial balance of zero.

As pointed out, equation (2.4) for the internal rate of return is a polynomial equation in  $c$  of degree  $n$ , which does not, in general, have an analytic solution. However, it is almost always easy to solve the equation with a computer. From algebraic theory it is known that such an equation always has at least one root, and may have as many as  $n$  roots, but some or all of these roots may be complex numbers. Fortunately the most common form of investment, where there is an initial cash outlay followed by several positive flows, leads to a unique positive solution. Hence the internal rate of return is then well defined and relatively easy to calculate. (See Exercise 4.) The formal statement of the existence of the positive root embodies the main result concerning the internal rate of return.

**Main theorem of internal rate of return** *Suppose the cash flow stream  $(x_0, x_1, \dots, x_n)$  has  $x_0 < 0$  and  $x_k \geq 0$  for all  $k$ ,  $k = 1, 2, \dots, n$ , with at least one term being strictly positive. Then there is a unique positive root to the equation*

$$0 = x_0 + x_1c + x_2c^2 + \cdots + x_nc^n.$$



**FIGURE 2.3 Function for proof.** If  $x_0 < 0$  and  $x_k \geq 0$  for all  $k$ ,  $1 \leq k \leq n$ , with at least one term being strictly positive, then the function  $f(c)$  will start below the horizontal axis and increase monotonically as  $c$  increases. Therefore there must be a unique positive solution  $c$  satisfying  $f(c) = 0$ .

Furthermore, if  $\sum_{k=0}^n x_k > 0$  (meaning that the total amount returned exceeds the initial investment), then the corresponding internal rate of return  $r = (1/c) - 1$  is positive.

**Proof:** We plot the function  $f(c) = x_0 + x_1c + x_2c^2 + \dots + x_nc^n$ , as shown in Figure 2.3. Note that  $f(0) < 0$ . However, as  $c$  increases, the value of  $f(c)$  also increases, since at least one of the cash flow terms is strictly positive. Indeed, it increases without limit as  $c$  increases to infinity. Since the function is continuous, it must cross the axis at some value of  $c$ . It cannot cross more than once, because it is strictly increasing. Hence there is a unique real value  $c_0$ , which is positive, for which  $f(c_0) = 0$ .

If  $\sum_{k=0}^n x_k > 0$ , which means that there is a net positive (nondiscounted) cash flow, then  $f(1) > 0$ . This means that the solution  $c_0$  satisfying  $f(c_0) = 0$  must be less than 1. Therefore  $r_0 = (1/c_0) - 1 > 0$ , where  $r_0$  is the internal rate of return. ■

If some (or all) solutions to the equation for internal rate of return are complex, the interpretation of these values is not simple. In general it is reasonable to select the solution that has the largest real part and use that real part to determine the internal rate of return. In practice, however, this is not often a serious issue, since suitable real roots typically exist.

## 2.5 Evaluation Criteria

The essence of investment is selection from a number of alternative cash flow streams. In order to do this intelligently, the alternative cash flow streams must be evaluated according to a logical and standard criterion. Several different criteria are used in

practice, but the two most important methods are those based on present value and on internal rate of return.

## Net Present Value

Present value evaluates alternatives by simply ranking them according to their present values—the higher the present value, the more desirable the alternative. When one uses present value this way, one must include *all* cash flows associated with the investment, both positive and negative. To emphasize that, the expression **net present value** (NPV) is frequently used. Net present value is the present value of the benefits minus the present value of the costs. Often, to emphasize this partition of benefits and costs, the terms **present worth of benefits** and **present worth of costs** are used, both of which are just present values. Net present value is the difference between these two terms. To be worthy of consideration, the cash flow stream associated with an investment must have a positive net present value.

**Example 2.4 (When to cut a tree)** Suppose that you have the opportunity to plant trees that later can be sold for lumber. This project requires an initial outlay of money in order to purchase and plant the seedlings. No other cash flow occurs until the trees are harvested. However, you have a choice as to when to harvest: after 1 year or after 2 years. If you harvest after 1 year, you get your return quickly; but if you wait an additional year, the trees will have additional growth and the revenue generated from the sale of the trees will be greater.

We assume that the cash flow streams associated with these two alternatives are

- (a)  $(-1, 2)$  cut early
- (b)  $(-1, 0, 3)$  cut later.

We also assume that the prevailing interest rate is 10%. Then the associated net present values are

- (a)  $\text{NPV} = -1 + 2/1.1 = .82$
- (b)  $\text{NPV} = -1 + 3/(1.1)^2 = 1.48.$

Hence according to the net present value criterion, it is best to cut later.

The net present value criterion is quite compelling, and indeed it is generally regarded as the single best measure of an investment's merit. It has the special advantage that the present values of different investments can be added together to obtain a meaningful composite. This is because the present value of a sum of cash flow streams is equal to the sum of the present values of the corresponding cash flows. Note, for example, that we were able to compare the two investment alternatives associated with tree farming even though the cash flows were at different times. In general, an

investor can compute the present value of individual investments and also the present value of an entire portfolio.

## Internal Rate of Return

Internal rate of return can also be used to rank alternative cash flow streams. The rule is simply this: the higher the internal rate of return, the more desirable the investment. However, a potential investment, or project, is presumably not worth considering unless its internal rate of return is greater than the prevailing interest rate. If the internal rate of return is greater than the prevailing interest rate, the investment is considered better than what is available externally in the financial market.

**Example 2.5 (When to cut a tree, continued)** Let us use the internal rate of return method to evaluate the two tree harvesting proposals considered in Example 2.4. The equations for the internal rate of return in the two cases are

- (a)  $-1 + 2c = 0$
- (b)  $-1 + 3c^2 = 0$ .

As usual,  $c = 1/(1+r)$ . These have the following solutions:

$$(a) c = \frac{1}{2} = \frac{1}{1+r}; \quad r = 1.0$$

$$(b) c = \frac{\sqrt{3}}{3} = \frac{1}{1+r}; \quad r = \sqrt{3} - 1 \approx .73.$$

In other words, for (a), cut early, the internal rate of return is 100%, whereas for (b) it is about 70%. Hence under the internal rate of return criterion, the best alternative is (a). Note that this is opposite to the conclusion obtained from the net present value criterion.

## Discussion of the Criteria

There is considerable debate as to which of the two criteria, net present value or internal rate of return, is the most appropriate for investment evaluation. Both have attractive features, and both have limitations. (As shown, they can even give conflicting recommendations.) Net present value is simplest to compute; it does not have the ambiguity associated with the several possible roots of the internal rate of return equation. Also net present value can be broken into component pieces, unlike internal rate of return. However, internal rate of return has the advantage that it depends only on the properties of the cash flow stream, and not on the prevailing interest rate (which in practice may not be easily defined). In fact, the two methods both have appropriate roles, but in different situations.

The primary difference between the two criteria can be explained in terms of the “when to cut a tree” example. We must look beyond the single cycle of tree farming to a series of cycles. Suppose that the proceeds of the first harvest are used to plant additional trees, starting a long series of expansion in the tree farming business. Under plan (a), cut early, the business can be doubled every year because the revenue received at the end of the year is twice that required at the beginning. In plan (b), cut later, the business can be tripled every 2 years by the same reasoning. Tripling every 2 years is equivalent, in the long run, to increasing by a factor of  $\sqrt{3} \approx 1.73$  every year. The yearly growth rates of these two plans, factors of 2 and  $\sqrt{3}$ , respectively, are each equal to 1 plus the internal rates of return of the plans—and this equality is true in general. So in this kind of situation, where the proceeds of the investment can be repeatedly reinvested in the same type of project but scaled in size, it makes sense to select the project with the largest internal rate of return—in order to get the greatest growth of capital.

On the other hand, suppose that this investment is a one-time opportunity and cannot be repeated. Here the net present value method is the appropriate criterion, since it compares the investment with what could be obtained through normal channels (which offer the prevailing rate of interest).

It is widely agreed (by theorists, but not necessarily by practitioners) that, overall, the best criterion is that based on net present value. If used intelligently, it will provide consistency and rationality, but a perspective broader than that of an isolated project often is required. In the case of cutting the trees, for example, an enlightened present value analysis will agree with the result obtained by the internal rate of return criterion. If the two possible futures are developed fully, corresponding to the two cutting policies, the present value criterion, applied to the long series of expanding cash flows, would also direct that plan (a) be adopted.

There are many other factors that influence a good present value analysis—and perhaps make such an analysis more complex than suggested by the direct formal statement of the criterion. One significant issue is the selection of the interest rate to be used in the calculation. In practice, there are several different “risk-free” rates of interest in the financial market: the rate paid by bank certificates of deposit, the 3-month U.S. Treasury bill rate, and the rate paid by the highest grade commercial bonds are examples. Furthermore, the rates for borrowing are typically slightly higher than those for lending. The difference between all these choices can be several percentage points. In business decisions it is common to use a figure called the **cost of capital** as the baseline rate. This figure is the rate of return that the company must offer to potential investors in the company; that is, it is the cost the company must pay to get additional funds. Or sometimes it is taken to be the rate of return expected on alternative desirable projects. However, some of these cost of capital figures are derived from uncertain cash flow streams and are not really appropriate measures of a risk-free interest rate. For present value calculations it is best to use rates that represent true interest rates, since we assume that the cash flows are certain. Some of the apparent differences in these rates are explained and justified in Chapter 4, but still there is room for judgment.

Another factor to consider is that present value by itself does not reveal much about the rate of return. Two alternative investments might each have a net present

value of \$100, but one might require an investment of \$100 whereas the other requires \$1 million. Clearly these two alternatives should be viewed differently. Net present value in its simplest form is not the whole story (but we never said it was). It forms a solid starting point, but one must supplement its use with additional structure.

## 2.6 Applications and Extensions\*

This section illustrates how the concepts of this chapter can be used to evaluate real investment opportunities and projects. Often creative thinking is required to capture the essence of a situation in a form that is suitable for analysis.

Not all of these applications need be read during the first pass through this chapter; but as one returns to the chapter, these examples should help clarify the underlying concepts.

### Net Flows

In conducting a cash flow analysis using either net present value or internal rate of return, it is essential that the net of income minus expense (that is, net profit) be used as the cash flow each period. The net profit usually can be found in a straightforward manner, but the process can be subtle in complex situations. In particular, taxes often introduce complexity because certain tax-accounting costs and profits are not always equal to actual cash outflows or inflows. Taxes are considered in a later subsection.

Here we use a relatively simple example involving a gold mine to illustrate net present value analysis. Various gold mine examples are used throughout the book to illustrate how, as we extend our conceptual understanding, we can develop deeper analyses of the same kind of investment. The Simplico gold mine example presented here is the simplest of the series.

**Example 2.6 (Simplico gold mine)** The Simplico gold mine has a great deal of remaining gold deposits, and you are part of a team that is considering leasing the mine from its owners for a period of 10 years. Gold can be extracted from this mine at a rate of up to 10,000 ounces per year at a cost of \$200 per ounce. This cost is the total operating cost of mining and refining, exclusive of the cost of the lease. Currently the market price of gold is \$400 per ounce. The interest rate is 10%. Assuming that the price of gold, the operating cost, and the interest rate remain constant over the 10-year period, what is the present value of the lease?

This is fairly straightforward. We ignore the lease expense and just find the present value of the operating profits. It is clear that the mine should be operated at full capacity every year, giving a profit of  $10,000 \times (\$400 - \$200) = \$2$  million per year. We assume that these cash flows occur at the end of each year.

---

\* Sections marked by stars may be skipped at first reading.

The cash flow stream therefore consists of 10 individual flows of \$2M (that is, \$2 million) at the end of each year. The present value is accordingly

$$PV = \sum_{k=1}^{10} \frac{\$2M}{(1.1)^k}.$$

This can be evaluated either by direct summation or by using the formula for the sum of a geometric series. (See Section 3.2.) The result is

$$PV = \$2M \left[ 1 - \left( \frac{1}{1.1} \right)^{10} \right] \times 10 = \$12.29M$$

and this is the value of the lease.

## Cycle Problems

When using interest rate theory to evaluate ongoing (repeatable) activities, it is essential that alternatives be compared over the same time horizon. The difficulties that can arise from not doing this are illustrated in the tree cutting example. The two alternatives in that example have different cycle lengths, but the nature of the possible repetition of the cycles was not clearly spelled out originally.

We illustrate here two ways to account properly for different cycle lengths. The first is to repeat each alternative until both terminate at the same time. For example, if a first alternative lasts 2 years and a second lasts 4 years, then two cycles of the first alternative are comparable to one of the second. The other method for comparing alternatives with different cycle lengths is to assume that an alternative will be repeated indefinitely. Then a simple equation can be written for the value of the entire infinite-length stream.

**Example 2.7 (Automobile purchase)** You are contemplating the purchase of an automobile and have narrowed the field down to two choices. Car A costs \$20,000, is expected to have a low maintenance cost of \$1,000 per year (payable at the beginning of each year after the first year), but has a useful mileage life that for you translates into 4 years. Car B costs \$30,000 and has an expected maintenance cost of \$2,000 per year (after the first year) and a useful life of 6 years. Neither car has a salvage value. The interest rate is 10%. Which car should you buy?

We analyze this choice by assuming that similar alternatives will be available in the future—we are ignoring inflation—so this purchase is one of a sequence of car purchases. To equalize the time horizon, we assume a planning period of 12 years, corresponding to three cycles of car A and two of car B.

We analyze simple cycles and combined cycles as follows.

Car A:

$$\begin{aligned} \text{One cycle} \quad PV_A &= 20,000 + 1,000 \sum_{k=1}^3 \frac{1}{(1.1)^k} \\ &= \$22,487 \end{aligned}$$

$$\begin{aligned} \text{Three cycles} \quad PV_{A3} &= PV_A \left[ 1 + \frac{1}{(1.1)^4} + \frac{1}{(1.1)^8} \right] \\ &= \$48,336 \end{aligned}$$

Car B:

$$\begin{aligned} \text{One cycle} \quad PV_B &= 30,000 + 2,000 \sum_{k=1}^5 \frac{1}{(1.1)^k} \\ &= \$37,582 \end{aligned}$$

$$\begin{aligned} \text{Two cycles} \quad PV_{B2} &= PV_B \left[ 1 + \frac{1}{(1.1)^6} \right] \\ &= \$58,795. \end{aligned}$$

Hence car A should be selected because its cost has the lower present value over the common time horizon.

**Example 2.8 (Machine replacement)** A specialized machine essential for a company's operations costs \$10,000 and has operating costs of \$2,000 the first year. The operating cost increases by \$1,000 each year thereafter. We assume that these operating costs occur at the end of each year. The interest rate is 10%. How long should the machine be kept until it is replaced by a new identical machine? Assume that due to its specialized nature the machine has no salvage value.

This is an example where the cash flow stream is not fixed in advance because of the unknown replacement time. We must also account for the cash flows of the replacement machines. This can be done by writing an equation having PV on both sides. For example, suppose that the machine is replaced every year. Then the cash flow (in thousands) is  $(-10, -2)$  followed by  $(0, -10, -2)$  and then  $(0, 0, -10, -2)$ , and so forth. However, we can write the total PV of the costs compactly as

$$PV = 10 + 2/1.1 + PV/1.1$$

because after the first machine is replaced, the stream from that point looks identical to the original one, except that this continuing stream starts 1 year later and hence must be discounted by the effect of 1 year's interest. The solution to this equation is  $PV = 130$  or, in our original units, \$130,000.

**TABLE 2.2**  
**MACHINE REPLACEMENT**

Replacement year	Present value
1	130,000
2	82,381
3	69,577
4	65,358
5	64,481
6	65,196

The total present value is found for various replacement frequencies. The best policy corresponds to the frequency having the smallest total present value.

We may do the same thing assuming 2-year replacement, then 3 years, and so forth. The general approach is based on the equation

$$PV_{\text{total}} = PV_1 \text{ cycle} + \left( \frac{1}{1.1} \right)^k PV_{\text{total}}$$

where  $k$  is the length of the basic cycle. This leads easily to Table 2.2.

From the table we see that the smallest present value of cost occurs when the machine is replaced after 5 years. Hence that is the best replacement policy.

## Taxes

Taxes can complicate a cash flow value analysis. No new conceptual issues arise; it is just that taxes can obscure the true definition of cash flow. If a uniform tax rate were applied to all revenues and expenses as taxes and credits, respectively, then recommendations from before-tax and after-tax analyses would be identical. The present value figures from the latter analysis would merely all be scaled by the same factor; that is, all would be multiplied by 1 minus the tax rate. The internal rate of return figures would be identical. Hence rankings using either net present value or internal rate of return would remain the same as those without taxes. For this reason taxes are ignored in many of our examples. Sometimes, however, the cash flows required to be reported to the government on tax forms are *not* true cash flows. This is why firms often must keep two sets of accounts—one for tax purposes and one for decision-making purposes. There is nothing illegal about this practice; it is a reality introduced by the tax code.

A tax-induced distortion of cash flows frequently accompanies the treatment of property depreciation. Depreciation is treated as a negative cash flow by the government, but the timing of these flows, as reported for tax purposes, rarely coincides with actual cash outlays.

Taxes can be effectively treated by use of the concept of an **after-tax rate**. To understand this, suppose you invest \$100 and the outside comparison rate is 10%.

This means that we expect the \$100 to grow to \$110 in one period. Suppose now that the gain is taxed at a rate of, say, 43%. This will leave you with  $\$100 + \$10 - \$4.3 = \$105.70$ . The present value of \$100 can be achieved if we discount our net after-tax amount at a discount factor of 5.70%. This is the rate that should be used for after-tax amounts. In general, if the outside comparison rate is  $r$  and the tax rate is  $t$ , the after-tax rate is  $(1 - t)r$ . In the example, this is  $(1 - .43)(.1) = .057$ .

**Example 2.9 (Depreciation)** Suppose a firm purchases a machine for \$10,000. This machine has a useful life of 4 years and its use generates a cash flow of \$3,000 each year. The machine has a salvage value of \$2,000 at the end of 4 years.

The government does not allow the full cost of the machine to be reported as an expense the first year, but instead requires that the cost of the machine be depreciated over its useful life. There are several depreciation methods, each applicable under various circumstances, but for simplicity we shall assume the straightline method. In this method a fixed portion of the cost is reported as depreciation each year. Hence corresponding to a 4-year life, one-fourth of the cost (minus the estimated salvage value) is reported as an expense deductible from revenue each year.

If we assume a combined federal and state tax rate of 43%, we obtain the cash flows, before and after tax, shown in Table 2.3. The present value at 10% of the original cash flow is \$876. The after-tax amounts assuming immediate application of the tax are shown in the second column. The present value of this using the after-tax rate is \$1,180. The result of the depreciation schedule is shown in the final column. This should be discounted at the after-tax rate of 5.7%, resulting in a present value of \$569. This much can be lost by the requirement to use depreciation.

## Inflation

Inflation is another factor that often causes confusion, arising from the choice between using actual dollar values to describe cash flows and using values expressed in

**TABLE 2.3**  
**CASH FLOWS BEFORE AND AFTER TAX**

Year	Before-tax cash flow	After-tax cash flow	Depreciation	Taxable income	Tax	After-tax cash flow
0	-10,000	-5,700				-10,000
1	3,000	1,710	2,000	1,000	430	2,570
2	3,000	1,710	2,000	1,000	430	2,570
3	3,000	1,710	2,000	1,000	430	2,570
4	5,000	2,850	2,000	1,000	430	4,570
PV @	10%	5.7%				5.7%
	876	1.180				569

*From a present value viewpoint, tax rules for treatment of depreciation can reduce the after-tax profitability of a project.*

purchasing power, determined by reducing inflated future dollar values back to a nominal level.

Inflation is characterized by an increase in general prices with time. Inflation can be described quantitatively in terms of an **inflation rate**  $f$ . Prices 1 year from now will on average be equal to today's prices multiplied by  $(1+f)$ . Inflation compounds much like interest does, so after  $k$  years of inflation at rate  $f$ , prices will be  $(1+f)^k$  times their original values. Of course, inflation rates do not remain constant, but in planning studies future rates are usually estimated as constant.

Another way to look at inflation is that it erodes the purchasing power of money. A dollar today does not purchase as much bread or milk, for example, as a dollar did 10 years ago. In other words, we can think of prices increasing or, alternatively, of the value of money decreasing. If the inflation rate is  $f$ , then the value of a dollar next year in terms of the purchasing power of today's dollar is  $1/(1+f)$ .

It is sometimes useful to think explicitly in terms of the same kind of dollars, eliminating the influence of inflation. Thus we consider **constant dollars** or, alternatively, **real dollars**, defined relative to a given reference year. These are the (hypothetical) dollars that continue to have the same purchasing power as dollars did in the reference year. These dollars are defined in contrast to the **actual** or **nominal dollars** that we really use in transactions.

This leads us to define a new interest rate, termed the **real interest rate**, which is the rate at which real dollars increase if left in a bank that pays the nominal rate. To understand the meaning of the real interest rate, imagine depositing money in the bank at time zero, then withdrawing it 1 year later. The purchasing power of the bank balance has probably increased in spite of inflation, and this increase measures the real rate of interest.

If one goes through that thinking, when  $r$  is the nominal interest rate and  $f$  is the inflation rate, it is easy to see that

$$1 + r_0 = \frac{1+r}{1+f},$$

where  $r_0$  denotes the real rate of interest. This equation expresses the fact that money in the bank increases (nominally) by  $1+r$ , but its purchasing power is deflated by  $1/(1+f)$ . We can solve for  $r_0$  as

$$r_0 = \frac{r-f}{1+f}. \quad (2.5)$$

Note that for small levels of inflation the real rate of interest is approximately equal to the nominal rate of interest minus the inflation rate.

A cash flow analysis can be carried out using either actual (nominal) dollars or real dollars, but the danger is that a mixture of the two might be used inadvertently. Such a mixture sometimes occurs in the planning studies in large corporations. The operating divisions, which are primarily concerned with physical inputs and outputs, may extrapolate real cash flows into the future. But corporate headquarters, being primarily concerned with the financial market and tax rules, may find the use of nominal (that is, actual) cash flows more convenient and hence may discount at the nominal rate. The result can be an undervaluation by headquarters of project proposals

**TABLE 2.4****INFLATION**

Year	Real cash flow	PV @5.77%	Nominal cash flow	PV @10%
0	-10,000	-10,000	-10,000	-10,000
1	5,000	4,727	5,200	4,727
2	5,000	4,469	5,408	4,469
3	5,000	4,226	5,624	4,226
4	3,000	2,397	3,510	2,397
Total		5,819		5,819

The projected real cash flows of the second column have the present values, at the real rate of interest, shown in the third column. The fourth column lists the cash flows that would occur under 4% inflation, and their present values at the 10% nominal rate of interest are given in the fifth column.

submitted by the divisions relative to valuations that would be obtained if inflation were treated consistently.

We illustrate now how an analysis can be carried out consistently by using either real or nominal cash flows.

**Example 2.10 (Inflation)** Suppose that inflation is 4%, the nominal interest rate is 10%, and we have a cash flow of real (or constant) dollars as shown in the second column of Table 2.4. (It is common to estimate cash flows in constant dollars, relative to the present, because “ordinary” price increases can then be neglected in a simple estimation of cash flows.) To determine the present value in real terms we must use the real rate of interest, which from (2.5) is  $r_0 = (.10 - .04)/1.04 = 5.77\%$ .

Alternatively, we may convert the cash flow to actual (nominal) terms by inflating the figures using the appropriate inflation factors. Then we determine the present value using the nominal interest rate of 10%. Both methods produce the same result.

## 2.7 Summary

The time value of money is expressed concretely as an interest rate. The 1-year interest rate is the price paid (expressed as a percentage of principal) for borrowing money for 1 year. In simple interest, the interest payment when borrowing money in subsequent years is identical in magnitude to that of the first year. Hence, for example, the bank balance resulting from a single deposit would grow linearly year by year. In compound interest, the interest payment in subsequent years is based on the balance at the beginning of that year. Hence the bank balance resulting from a single deposit would grow geometrically year by year.

A useful approximation is that the number of years required for a deposit to double in value when compounded yearly is  $72/i$ , where  $i$  is the interest rate expressed as a percentage. For example, at 10%, money doubles in about 7 years.

Interest can be compounded at any frequency, not just yearly. It is even possible to compound continuously, which leads to bank balances that grow exponentially with

time. When interest is compounded more frequently than yearly, it is useful to define both a nominal rate and an effective annual rate of interest. The nominal rate is the rate used for a single period divided by the length (in years) of a period. The effective rate is the rate that, if applied without compounding, would give the same total balance for money deposited for one full year. The effective rate is larger than the nominal rate. For example, an 8% nominal annual rate corresponds to an 8.24% effective annual rate under quarterly compounding.

Money received in the future is worth less than the same amount of money received in the present because money received in the present can be loaned out to earn interest. Money to be received at a future date must be discounted by dividing its magnitude by the factor by which present money would grow if loaned out to that future date. There is, accordingly, a discount factor for each future date.

The present value of a cash flow stream is the sum of the discounted magnitudes of the individual cash flows of the stream. An ideal bank can transform a cash flow stream into any other with the same present value.

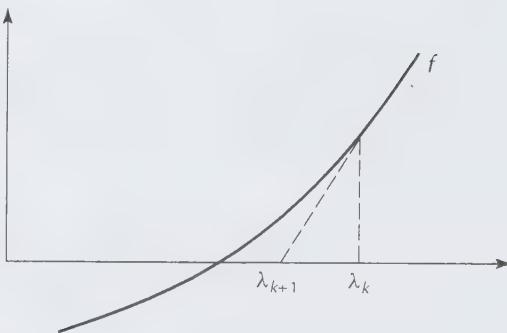
The internal rate of return of a cash flow stream is an interest rate that, if used to evaluate the present value of the stream, would cause that present value to be zero. In general, this rate is not well defined. However, when the cash flow stream has an initial negative flow followed by positive flows, the internal rate of return is well defined.

Present value and internal rate of return are the two main methods used to evaluate proposed investment projects that generate deterministic cash flow streams. Under the present value framework, if there are several competing alternatives, then the one with the highest present value should be selected. Under the internal rate of return criterion, the alternative with the largest internal rate of return should be selected.

Analyses using these methods are not always straightforward. In particular, consideration of various cycle lengths, taxes, and inflation each require careful attention.

## Exercises

1. (A nice inheritance) Suppose \$1 were invested in 1776 at 3.3% interest compounded yearly.
  - (a) Approximately how much would that investment be worth today: \$1,000, \$10,000, \$100,000, or \$1,000,000?
  - (b) What if the interest rate were 6.6%?
2. (The 72 rule) The number of years  $n$  required for an investment at interest rate  $r$  to double in value must satisfy  $(1+r)^n = 2$ . Using  $\ln 2 = .69$  and the approximation  $\ln(1+r) \approx r$  valid for small  $r$ , show that  $n \approx 69/i$ , where  $i$  is the interest rate percentage (that is,  $i = 100r$ ). Using the better approximation  $\ln(1+r) \approx r - \frac{1}{2}r^2$ , show that for  $r \approx .08$  there holds  $n \approx 72/i$ .
3. (Effective rates) Find the corresponding effective rates for:
  - (a) 3% compounded monthly.
  - (b) 18% compounded monthly.
  - (c) 18% compounded quarterly.

**FIGURE 2.4** Newton's method.

4. (Newton's method ◊) The IRR is generally calculated using an iterative procedure. Suppose that we define  $f(\lambda) = -a_0 + a_1\lambda + a_2\lambda^2 + \dots + a_n\lambda^n$ , where all  $a_i$ 's are positive and  $n > 1$ . Here is an iterative technique that generates a sequence  $\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_k, \dots$  of estimates that converges to the root  $\bar{\lambda} > 0$ , solving  $f(\bar{\lambda}) = 0$ . Start with any  $\lambda_0 > 0$  close to the solution. Assuming  $\lambda_k$  has been calculated, evaluate

$$f'(\lambda_k) = a_1 + 2a_2\lambda_k + 3a_3\lambda_k^2 + \dots + na_n\lambda_k^{n-1}$$

and define

$$\lambda_{k+1} = \lambda_k - \frac{f(\lambda_k)}{f'(\lambda_k)}.$$

This is Newton's method. It is based on approximating the function  $f$  by a line tangent to its graph at  $\lambda_k$ , as shown in Figure 2.4. Try the procedure on  $f(\lambda) = -1 + \lambda + \lambda^2$ . Start with  $\lambda_0 = 1$  and compute four additional estimates.

5. (Tree cut) Suppose that you have the opportunity to plant trees that later can be sold for lumber. This project requires an initial outlay of money in order to purchase and plant the seedlings. No other cash flow occurs until the trees are harvested. However, you have a choice as to when to harvest. If you harvest after 1 year, you get your return quickly; but if you wait, the trees will have additional growth and the revenue generated from the sale of the trees will be greater. Assume that the cash flow streams associated with these alternatives are:

Wait 1 year:  $(-1, 2)$

Wait 2 years:  $(-1, 0, 3)$

Wait 3 years:  $(-1, 0, 0, 4)$

⋮

Wait  $n$  years:  $(-1, 0, 0, \dots, 0, 0, n+1)$

The prevailing interest rate is 10%. When is the best time to cut the trees?

6. (Equivalent rates) What rate of interest (to five digits) is equivalent to 10% yearly under
- monthly compounding?
  - continuous compounding?

---

◊ Exercises followed by ◊ are mathematically more difficult than average.

7. (A prize) A major lottery advertises that it pays the winner \$10 million. However, this prize money is paid at the rate of \$500,000 each year (with the first payment being immediate) for a total of 20 payments. What is the present value of this prize at 10% interest?
8. (Sunk costs) A young couple has made a nonrefundable deposit of the first month's rent (equal to \$1,000) on a 6-month apartment lease. The next day they find a different apartment that they like just as well, but its monthly rent is only \$900. They plan to be in the apartment only 6 months. Should they switch to the new apartment? What if they plan to stay 1 year? Assume an interest rate of 12%.
9. (Shortcut) Gavin Jones is inquisitive and determined to learn both the theory and the application of investment theory. He pressed the tree farmer for additional information and learned that it was possible to delay cutting the trees of Example 2.4 for another year. The farmer said that, from a present value perspective, it was not worthwhile to do so. Gavin instantly deduced that the revenue obtained must be less than  $x$ . What is  $x$ ?
10. (Copy machines  $\oplus$ ) Two copy machines are available. Both have useful lives of 5 years. One machine can be either leased or purchased outright; the other must be purchased. Hence there are a total of three options: A, B, and C. The details are shown in Table 2.5. (The first year's maintenance is included in the initial cost. There are then four additional maintenance payments, occurring at the beginning of each year, followed by revenues from resale.) The present values of the expenses of these three options using a 10% interest rate are also indicated in the table. According to a present value analysis, the machine of least cost, as measured by the present value, should be selected; that is, option B.

**TABLE 2.5**  
**COPY MACHINE OPTIONS**

	Option		
	A	B	C
Initial outlay	6,000	30,000	35,000
Yearly expense	8,000	2,000	1,600
Resale value	0	10,000	12,000
Present value (@ 10%)	31,359	30,131	32,621

*Option A is a lease; options B and C are purchases of two alternative machines. All have 5-year lives.*

It is not possible to compute the IRR for any of these alternatives, because all cash flows are negative (except for the resale values). However, it is possible to calculate the IRR on an incremental basis. Find the IRR corresponding to a change from A to B. Is a change from A to B justified on the basis of the IRR?

11. (An appraisal) You are considering the purchase of a nice home. It is in every way perfect for you and in excellent condition, except for the roof. The roof has only 5 years of life remaining. A new roof would last 20 years, but would cost \$20,000. The house is expected to last forever. Assuming that costs will remain constant and that the interest rate is 5%, what value would you assign to the existing roof?

---

$\oplus$  Exercises followed by  $\oplus$  require numerical computation.

- 12.** (Oil depletion allowance  $\oplus$ ) A wealthy investor spends \$1 million to drill and develop an oil well that has estimated reserves of 200,000 barrels. The well is to be operated over 5 years, producing the estimated quantities shown in the second column of Table 2.6. It is estimated that the oil will be sold for \$20 per barrel. The net income is also shown.

**TABLE 2.6**  
**OIL INVESTMENT DETAILS**

Year	Barrels produced	Gross revenue	Net income	Option 1	Option 2	Depletion allowance	Taxable income	Tax	After-tax income
1	80,000	1,600,000	1,200,000	352,000	400,000	400,000	800,000	360,000	840,000
2	70,000	1,400,000	1,000,000						
3	50,000	1,000,000	500,000						
4	30,000	600,000	200,000						
5	10,000	200,000	50,000						

A depletion allowance, for tax purposes, can be computed in either of two ways each year: 22% of gross revenue up to 50% of net income before such deduction (option 1), or the investment cost of the product, equal in this case to the unit cost of the reserves, \$5 per barrel (option 2). The allowance is deducted from the net income to determine the taxable income. The investor is in the 45% tax bracket.

- (a) Complete Table 2.6 and show that the total depletion allowance exceeds the original investment.
- (b) Calculate the PV and the IRR for this investment. Assume an interest rate of 20%.

- 13.** (Conflicting recommendations  $\oplus$ ) Consider the two projects whose cash flows are shown in Table 2.7. Find the IRRs of the two projects and the NPVs at 5%. Show that the IRR and NPV figures yield different recommendations. Can you explain this?

**TABLE 2.7**

	Years					
	0	1	2	3	4	5
Project 1	-100	30	30	30	30	30
Project 2	-150	42	42	42	42	42

- 14.** (Domination) Suppose two competing projects have cash flows of the form  $(-A_1, B_1, B_1, \dots, B_1)$  and  $(-A_2, B_2, B_2, \dots, B_2)$ , both with the same length and  $A_1, A_2, B_1, B_2$  all positive. Suppose  $B_1/A_1 > B_2/A_2$ . Show that project 1 will have a higher IRR than project 2.

- 15.** (Crossing  $\diamond$ ) In general, we say that two projects with cash flows  $x_i$  and  $y_i$ ,  $i = 0, 1, 2, \dots, n$ , cross if  $x_0 < y_0$  and  $\sum_{i=0}^n x_j > \sum_{i=0}^n y_i$ . Let  $P_x(d)$  and  $P_y(d)$  denote the present values of these two projects when the discount factor is  $d$ .

- (a) Show that there is a crossover value  $c > 0$  such that  $P_x(c) = P_y(c)$ .  
 (b) For Exercise 13, calculate the crossover value  $c$ .
- 16.** (Depreciation choice) In the United States the accelerated cost recovery system (ACRS) must be used for depreciation of assets placed into service after December 1980. In this system, assets are classified into categories specifying the effective tax life. The classification of “3-year property,” for example, includes automobiles, tractors for hauling highway trailers, light trucks, and certain manufacturing tools. The percentages of the cost for 3-year property that can be deducted for each of the first 3 years after purchase (including the year of purchase) are 25%, 38%, and 37%, respectively. The tax code also allows the alternate ACRS method, which for 3-year property means that the straight-line percentage of  $33\frac{1}{3}\%$  can be used for 3 years.  
 Which of these methods is preferred by an individual who wishes to maximize the present value of depreciation? How does the choice depend on the assumed rate of interest?
- 17.** (An erroneous analysis) A division of ABBOX Corporation has developed the concept of a new product. Production of the product would require \$10 million in initial capital expenditure. It is anticipated that 1 million units would be sold each year for 5 years, and then the product would be obsolete and production would cease. Each year’s production would require 10,000 hours of labor and 100 tons of raw material. Currently the average wage rate is \$30 per hour and the cost of the raw material is \$ 100 per ton. The product would sell for \$3.30 per unit, and this price is expected to be maintained (in real terms). ABBOX management likes to use a 12% discount rate for projects of this type and faces a 34% tax rate on profit. The initial capital expenditure can be depreciated in a straight line fashion over 5 years. In its first analysis of this project, management did not apply inflation factors to the extrapolated revenues and operating costs. What present value did they obtain? How would the answer change if an inflation rate of 4% were applied?

## References

The theory of interest, compounding, present value, and internal rate of return is covered extensively in many excellent textbooks. A few investment-oriented texts which discuss general notions of interest are [1–3]. The use of the concepts of NPV and IRR for ranking investment alternatives is developed in detail in the field of engineering economy. Excellent texts in that field include [4–5]. A more advanced study of interest is [6], which contains a continuous-time version of the “when to cut a tree” example, which inspired the example given in Section 2.5. Exercise 12 is a modification of an example in [4]. See [7] for an approach to the IRR multiple roots issue.

1. Bodie, H. M., A. Kane, and A. J. Marcus (2004), *Investments*, 6th ed., Irwin, Homewood, IL.
2. Brealey, R., and S. Meyers (2010), *Principles of Corporate Finance*, 10th ed., McGraw-Hill, New York.
3. Haugen, R. A. (2000), *Modern Investment Theory*, 5th ed., Prentice Hall, Englewood Cliffs, NJ.
4. Sullivan, W.G., E. M. Wicks, and C. P. Koelling (2011), *Engineering Economy*, 17th ed., Prentice-Hall, Englewood Cliffs, NJ.
5. Newman, D. G., T. G. Eschenbach, and J. P. Lavelle (2011), *Engineering Economic Analysis*, 11th ed., Oxford University Press, New York.
6. Hirshleifer, J. (1970), *Investment, Interest, and Capital*, Prentice Hall, Englewood Cliffs.
7. Hazen, G. B. (2003), “A New Perspective on Multiple Internal Rates of Return,” *The Engineering Economist*, **48**, 31–51.

# 3

## FIXED-INCOME SECURITIES

**A**n interest rate is a price, or rent, for the most popular of all traded commodities—money. The one-year interest rate, for example, is just the price that must be paid for borrowing money for one year. Markets for money are well developed, and the corresponding basic market price—interest—is monitored by everyone who has a serious concern about financial activity.

As shown in the previous chapter, the market interest rate provides a ready comparison for investment alternatives that produce cash flows. This comparison can be used to evaluate any cash flow stream: whether arising from transactions between individuals, associated with business projects, or generated by investments in securities.

However, the overall market associated with interest rates is more complex than the simple bank accounts discussed in the last chapter. Vast assortments of bills, notes, bonds, annuities, futures contracts, and mortgages are part of the well-developed markets for money. These market items are not real goods (or hard assets) in the sense of having intrinsic value—such as potatoes or gold—but instead are traded only as pieces of paper, or as entries in a computer database. These items, in general, are referred to as **financial instruments**. Their values are derived from the promises they represent. If there is a well-developed market for an instrument, so that it can be traded freely and easily, then that instrument is termed a **security**. There are many financial instruments and securities that are directly related to interest rates and, therefore, provide access to income—at a price defined by the appropriate interest rate or rates.

**Fixed-income securities** are financial instruments that are traded in well-developed markets and promise a fixed (that is, definite) income to the holder over a

span of time. In our terminology, they represent the ownership of a definite cash flow stream.

Fixed-income securities are important to an investor because they define the market for money, and most investors participate in this market. These securities are also important as additional comparison points when conducting analyses of investment opportunities that are not traded in markets, such as a firm's research projects, oil leases, and royalty rights. A comprehensive study of financial instruments most naturally starts with a study of fixed-income securities.

## 3.1 The Market for Future Cash

The classification of a security as being a fixed-income security is actually a bit vague. Originally this classification meant, as previously stated, that the security pays a fixed, well-defined cash flow stream to the owner. The only uncertainties about the promised stream were associated with whether the issuer of the security might **default** (by, say, going bankrupt), in which case the income would be discontinued or delayed. Now, however, some "fixed-income" securities promise cash flows whose magnitudes are tied to various contingencies or fluctuating indices. For example, payment levels on an adjustable-rate mortgage may be tied to an interest rate index, or corporate bond payments may in part be governed by a stock price. But in common parlance, such variations are allowed within a broader definition of fixed-income securities. The general idea is that a fixed-income security has a cash flow stream that is fixed except for variations due to well-defined contingent circumstances.

There are many different kinds of fixed-income securities, and we cannot provide a comprehensive survey of them here. However, we shall mention some of the principal types of fixed-income securities in order to indicate the general scope of such securities.

### Savings Deposits

Probably the most familiar fixed-income instrument is an interest-bearing bank deposit. These are offered by commercial banks, savings and loan institutions, and credit unions. In the United States most such deposits are guaranteed by agencies of the federal government. The simplest **demand deposit** pays a rate of interest that varies with market conditions. Over an extended period of time, such a deposit is not strictly of a fixed-income type; nevertheless, we place it in the fixed-income category. The interest is guaranteed in a **time deposit account**, where the deposit must be maintained for a given length of time (such as 6 months), or else a penalty for early withdrawal is assessed. A similar instrument is a **certificate of deposit** (CD), which is issued in standard denominations such as \$10,000. Large-denomination CDs can be sold in a market, and hence they qualify as securities.

## Money Market Instruments

The term **money market** refers to the market for short-term (1 year or less) loans by corporations and financial intermediaries, including, for example, banks. It is a well-organized market designed for large amounts of money, but it is not of great importance to long-term investors because of its short-term and specialized nature. Within this market **commercial paper** is the term used to describe unsecured loans (that is, loans without collateral) to corporations. The larger denominations of CDs mentioned earlier are also part of this market.

A **banker's acceptance** is a more involved money market instrument. If company A sells goods to company B, company B might send a written promise to company A that it will pay for the goods within a fixed time, such as 3 months. Some bank *accepts* the promise by promising to pay the bill on behalf of company B. Company A can then sell the banker's acceptance to someone else at a discount before the time has expired.

**Eurodollar deposits** are deposits denominated in dollars but held in a bank outside the United States. Likewise **Eurodollar CDs** are CDs denominated in dollars and issued by banks outside the United States. A distinction between these Eurodollars and regular dollars is due to differences in banking regulations and insurance.

## U.S. Government Securities

The U.S. Government obtains loans by issuing various types of fixed-income securities. These securities are usually considered to be of the highest credit quality since they are backed by the government itself. The most important government securities are sketched here.

**U.S. Treasury bills** are issued in denominations of \$10,000 or more with fixed terms to maturity of 13, 26, and 52 weeks. They are sold on a discount basis. Thus a bill with a face value of \$ 10,000 may sell for \$9,500, the difference between the price and the face value providing the interest. A bill can be redeemed for the full face value at the maturity date. New bills are offered each week and are sold at auction. They are highly **liquid** (that is, there is a ready market for them); hence they can be easily sold prior to the maturity date.

**U.S. Treasury notes** have maturities of 1 to 10 years and are sold in denominations as small as \$1,000. The owner of such a note receives a **coupon payment** every 6 months until maturity. This coupon payment represents an interest payment and its magnitude is fixed throughout the life of the note. At maturity the note holder receives the last coupon payment and the face value of the note. Like Treasury bills, these notes are sold at auction.

**U.S. Treasury bonds** are issued with maturities of more than 10 years. They are similar to Treasury notes in that they make coupon payments.

**U.S. Treasury inflation-protected securities (TIPS)** are adjusted for inflation by changing the principal according to the Consumer Price Index. The coupon rate remains constant, but the corresponding interest payment varies with the adjusted principal value.

**U.S. Treasury strips** are Treasury bonds sold in the secondary market in stripped form. Here each of the coupons is issued separately, as is the principal. So a 10-year bond when stripped will consist of 20 semiannual coupon securities (each with a separate CUSIP<sup>1</sup>) and an additional principal security. Each of these securities generates a single cash flow, with no intermediate coupon payments. Such a security is termed a **zero-coupon bond**.

## Other Bonds

Bonds are issued by agencies of the federal government, by state and local governments, and by corporations.

**Municipal bonds** are issued by agencies of state and local governments. There are two main types: **general obligation bonds**, which are backed by a governing body such as the state; and **revenue bonds**, which are backed either by the revenue to be generated by the project that will initially be funded by the bond issue or by the agency responsible for the project.

The interest income associated with municipal bonds is exempt from federal income tax and from state and local taxes in the issuing state. This feature means that investors are willing to accept lower interest rates on these bonds compared to other securities of similar quality.

**Corporate bonds** are issued by corporations for the purpose of raising capital for operations and new ventures. They vary in quality depending on the strength of the issuing corporation and on certain features of the bond itself.

Some corporate bonds are traded on an exchange, but most are traded over-the-counter in a network of bond dealers. These over-the-counter bonds are less liquid in the sense that there may be only a few trades per day of a particular issue.

A bond carries with it an **indenture**, which is a contract of terms. Some features that might be included are:

**Callable bonds** A bond is callable if the issuer has the right to repurchase the bond at a specified price. Usually this call price falls with time, and often there is an initial call protection period wherein the bond cannot be called.

**Sinking funds** Rather than incur the obligation to pay the entire face value of a bond issue at maturity, the issuer may establish a sinking fund to spread this obligation out over time. Under such an arrangement the issuer may repurchase a certain fraction of the outstanding bonds each year at a specified price.

**Debt subordination** To protect bond holders, limits may be set on the amount of additional borrowing by the issuer. Also the bondholders may be guaranteed that in the event of bankruptcy, payment to them takes priority over payments of other debt—the other debt being subordinated.

---

<sup>1</sup> The Committee on Uniform Securities Identification Procedures (CUSIP) assigns identifying CUSIP numbers and codes to all securities.

## Mortgages

To a typical homeowner, a mortgage looks like the opposite of a bond. A future homeowner usually will *sell* a home mortgage to generate immediate cash to pay for a home, obligating him- or herself to make periodic payments to the mortgage holder. The standard mortgage is structured so that equal monthly payments are made throughout its term, which contrasts to most bonds, which have a final payment equal to the face value at maturity. Most standard mortgages allow for early repayment of the balance. Hence from the mortgage holder's viewpoint the income stream generated is not completely fixed, since it may be terminated with an appropriate lump-sum payment at the discretion of the homeowner.

There are many variations on the standard mortgage. There may be modest-sized periodic payments for several years followed by a final **balloon payment** that completes the contract. **Adjustable-rate mortgages** adjust the effective interest rate periodically according to an interest rate index, and hence these mortgages do not really generate fixed income in the strict sense.

Mortgages are not usually thought of as securities, since they are written as contracts between two parties, for example, a homeowner and a bank. However, mortgages are typically "bundled" into large packages and traded among financial institutions. These **mortgage-backed securities** are quite liquid.

## Annuities

An **annuity** is a contract that pays the holder (the **annuitant**) money periodically, according to a predetermined schedule or formula, over a period of time. Pension benefits often take the form of annuities. Sometimes annuities are structured to provide a fixed payment every year for as long as the annuitant is alive, in which case the price of the annuity is based on the age of the annuitant when the annuity is purchased and on the number of years until payments are initiated.

There are numerous variations. Sometimes the level of the annuity payments is tied to the earnings of a large pool of funds from which the annuity is paid, sometimes the payments vary with time, and so forth.

Annuities are not really securities, since they are not traded. (The issuer certainly would not allow a change in annuitant if payments are tied to the life of the owner; likewise, an annuitant would not allow the annuity company to transfer their obligation to another company which might be less solvent.) Annuities are, however, considered to be investment opportunities that are available at standardized rates. Hence from an investor's viewpoint, they serve the same role as other fixed-income instruments.

## 3.2 Value Formulas

Many fixed-income instruments include an obligation to pay a stream of equal periodic cash flows. This is characteristic of standard coupon bonds that pay the holder a fixed sum on a regular basis; it also is characteristic of standard mortgages, of many

annuities, of standard automobile loans, and of other consumer loans. It is therefore useful to recognize that the present value of such a constant stream can be determined by a compact formula. This formula is difficult to evaluate by hand, and hence professionals working each day with such financial instruments typically have available appropriate tables, handheld calculators, or computer programs that relate present value to the magnitude and term of periodic payments. There are, for example, extensive sets of mortgage tables, bond tables, annuity rate tables, and so forth. We shall develop the basic formula here and illustrate its use.

## Perpetual Annuities

As a step toward the development of the formula we consider an interesting and conceptually useful fixed-income instrument termed a **perpetual annuity**, or **perpetuity**, which pays a fixed sum periodically *forever*. For example, it might pay \$1,000 every January 1 forever. Such annuities are quite rare (although such instruments actually do exist in Great Britain, where they are called **consols**).

The present value of a perpetual annuity can be easily derived. Suppose an amount  $A$  is paid at the end of each period, starting at the end of the first period, and suppose the *per-period* interest rate is  $r$ . Then the present value is

$$P = \sum_{k=1}^{\infty} \frac{A}{(1+r)^k}.$$

The terms in the summand represent a geometric series, and this series can be summed easily using a standard formula. Alternatively, if you have forgotten the standard formula, we can derive it by noting that

$$P = \sum_{k=1}^{\infty} \frac{A}{(1+r)^k} = \frac{A}{1+r} + \sum_{k=2}^{\infty} \frac{A}{(1+r)^k} = \frac{A}{1+r} + \frac{P}{1+r}.$$

We can solve this equation to find  $P = A/r$ . Hence we have the following basic result:

**Perpetual annuity formula** *The present value  $P$  of a perpetual annuity that pays an amount  $A$  every period, beginning one period from the present, is*

$$P = \frac{A}{r},$$

*where  $r$  is the one-period interest rate.*

**Example 3.1 (Perpetual annuity)** Consider a perpetual annuity of \$1,000 each year. At 10% interest its present value is

$$P = \frac{1,000}{.10} = \$10,000.$$

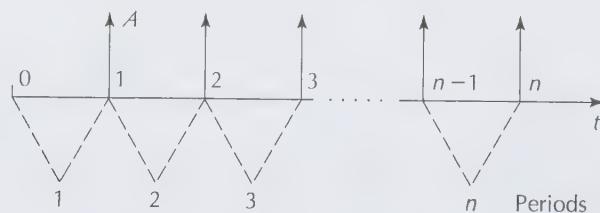
## Finite-Life Streams

Of more practical importance is the case where the payment stream has a finite lifetime. Suppose that the stream consists of  $n$  periodic payments of amount  $A$ , starting at the end of the current period and ending at period  $n$ . The pattern of periodic cash flows together with the time indexing system is shown in Figure 3.1.

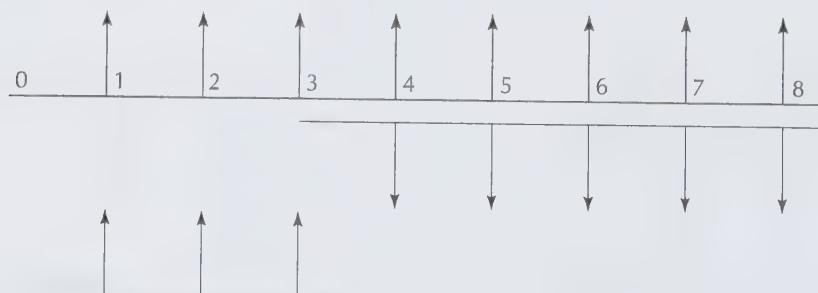
The present value of the finite stream relative to the interest rate  $r$  per period is

$$P = \sum_{k=1}^n \frac{A}{(1+r)^k}.$$

This is the sum of a finite geometric series. If you do not recall the formula for this sum, we can derive it easily by a simple trick. The value can be found by considering two perpetual annuities. Both pay an amount  $A$  each year, but one starts at time 1 and the other starts at time  $n+1$ . We *subtract* the second from the first. The result is the same as the original stream of finite life. This combination is illustrated in Figure 3.2 for the case of a stream of length 3.



**FIGURE 3.1 Time indexing.** Time is indexed from 0 to  $n$ . A period is a span between time points, with the first period being the time from 0 to 1. A standard annuity has a constant cash flow at the end of each period.



**FIGURE 3.2 Finite stream from two perpetual annuities.** The top line shows a perpetuity starting at time 1, the second a negative perpetuity starting at time 4. The sum of these two is a finite-life annuity with payments starting at time 1 and ending at time 3.

The value of the delayed annuity is found by discounting that annuity by the factor  $(1+r)^{-n}$  because it is delayed  $n$  periods. Hence we may write

$$P = \frac{A}{r} - \frac{A}{r(1+r)^n} = \frac{A}{r} \left[ 1 - \frac{1}{(1+r)^n} \right].$$

We now highlight this important result:

**Annuity formulas** Consider an annuity that begins payment one period from the present, paying an amount  $A$  each period for a total of  $n$  periods. The present value  $P$ , the one-period annuity amount  $A$ , the one-period interest rate  $r$ , and the number of periods  $n$  of the annuity are related by

$$P = \frac{A}{r} \left[ 1 - \frac{1}{(1+r)^n} \right]$$

or, equivalently,

$$A = \frac{r(1+r)^n P}{(1+r)^n - 1}.$$

Although these formulas are simple in concept and quite easy to derive, they are sufficiently complex that they cannot be evaluated easily by hand. It is for this reason that financial tables and financial calculators are commonly available. Professional tables of this type occupy several pages and typically give  $P/A$  as a function of  $r$  and  $n$ . For some purposes  $A/P$  (just the reciprocal) is more convenient, and there are tables written both ways.

It is important to note that in the formulas of this section,  $r$  is expressed as a per-period interest rate. If the period length is not equal to 1 year, this  $r$  will *not* be equal to the yearly rate; so care must be exercised.

The annuity formula is frequently used in the reverse direction; that is,  $A$  as a function of  $P$ . This determines the periodic payment that is equivalent (under the assumed interest rate) to an initial payment of  $P$ . This process of substituting periodic payments for a current obligation is referred to as **amortization**. Hence one may **amortize** the cost of an automobile over 5 years by taking out a 5-year loan.

**Example 3.2 (Loan calculation)** Suppose you have borrowed \$1,000 from a credit union. The terms of the loan are that the yearly interest is 12% compounded monthly. You are to make equal monthly payments of such magnitude as to repay (amortize) this loan over 5 years. How much are the monthly payments?

Five years is 60 months, and 12% a year compounded monthly is 1% per month. Hence we use the formula for  $n = 60$ ,  $r = 1\%$ , and  $P = \$1,000$ . We find that the payments  $A$  are \$22.20 per month.

**Example 3.3 (APR)** A typical advertisement from a mortgage broker is shown in Table 3.1. In addition to the interest rate, term of the loan, and maximum amount, there are listed points and the annual percentage rate (APR), which describe fees and expenses. **Points** is the percentage of the loan amount that is charged for providing the mortgage. Typically, there are additional expenses as well. All of these fees and

**TABLE 3.1**  
**MORTGAGE BROKER ADVERTISEMENT**

Rate	Pts	Term	Max amt	APR
7.625	1.00	30 yr	\$203,150	7.883
7.875	.50	30 yr	\$203,150	8.083
8.125	2.25	30 yr	\$600,000	8.399
7.000	1.00	15 yr	\$203,150	7.429
7.500	1.00	15 yr	\$600,000	7.859

Call 555-1213

Real Estate Broker, CA Dept. of Real Estate,  
Mortgage Masters, Inc.  
Current Fixed Rates

*APR is the rate of interest that implicitly includes the fees associated with a mortgage.*

expenses are added to the loan balance, and the sum is amortized at the stated rate over the stated period. This results in a fixed monthly payment amount  $A$ .

The **APR** is the rate of interest that, if applied to the loan amount without fees and expenses, would result in a monthly payment of  $A$ , exactly as before.

As a concrete example, suppose you took out a mortgage corresponding to the first listing in Table 3.1. Let us calculate the total fees and expenses. Using the APR of 7.883%, a loan amount of \$203,150, and a 30-year term, we find a monthly payment of  $A = \$1,474$ .

Now using an interest rate of 7.625% and the monthly payment calculated, we find a total initial balance of \$208,267. The total of fees and expenses is therefore  $\$208,267 - \$203,150 = \$5,117$ . The loan fee itself is 1 point, or \$2,032. Hence other expenses are  $\$5,117 - \$2,032 = \$3,085$ .

## Running Amortization\*

The formulas for amortization can be looked at in another way, linked directly to common accounting practice. Consider the loan of \$1,000 discussed in Example 3.2, which you will repay over 5 years at 12% interest (compounded monthly). Suppose you took out the loan on January 1, and the first payment is due February 1. The repayment process can be viewed as credits to a running monthly account. The account has an initial balance equal to the value of the loan—the original principal. Each month this balance is increased by an interest charge of 1% and then reduced by the payment amount. Assuming that you make payments as scheduled, the balance will decrease each month, reaching zero after 60 months. On July 1 you might receive a 6-month accounting such as that shown in Table 3.2, which illustrates how the balance decreases as payments are made.

It is common to regard each payment as consisting of two parts. The first part is the current interest; the second is a partial repayment of the principal. The running balance account procedure is consistent with reamortizing the loan each

**TABLE 3.2**  
**STATEMENT OF ACCOUNT TRANSACTIONS**

	Previous balance	Current interest	Payment received	New balance
January 1				1,000.00
February 1	1,000.00	10.00	22.20	987.80
March 1	987.80	9.88	22.20	975.48
April 1	975.48	9.75	22.20	963.03
May 1	963.03	9.63	22.20	950.46
June 1	950.46	9.50	22.20	937.76

*Each month the previous balance accumulates interest and is reduced by the current payment. The balance will be zero at the end of the loan term.*

month. Specifically, assuming all payments to date were made on schedule and of the proper amount, the payment level predicted by the formula to amortize the current balance over the months remaining in the original contract will always be \$22.20. For example, based on the July statement, one can amortize the balance of \$ 937.76 at 12% on June 1 (after making the June 1 payment) over a period of 55 months. The monthly payment required by this amortization would be \$22.20.

## Annual Worth\*

The annuity framework provides an alternative method for expressing a net present value analysis. This **annual worth** method has the advantage that it expresses its results in terms of a constant level of cash flow and thus is easily understood.

Suppose a project has an associated cash flow stream  $(x_0, x_1, \dots, x_n)$  over  $n$  years. A present value analysis uses a (fictitious) constant ideal bank with interest rate  $r$  to transform this stream hypothetically into an equivalent one of the form  $(v, 0, 0, \dots, 0)$ , where  $v$  is the net present value of the stream.

An annual worth analysis uses the same ideal bank to hypothetically transform the sequence to one of the form  $(0, A, A, A, \dots, A)$ . The value  $A$  is the annual worth (over  $n$  years) of the project. It is the equivalent net amount that is generated by the project if all amounts are converted to a fixed  $n$ -year annuity starting the first year.

Clearly  $A > 0$  exactly when  $v > 0$ , so the condition for acceptance of the project based on whether  $A > 0$  coincides with the net present value criterion.

**Example 3.4 (A capital cost)** The purchase of a new machine for \$100,000 (at time zero) is expected to generate additional revenues of \$25,000 for the next 10 years starting at year 1. If the discount rate is 16%, is this a profitable investment?

We simply need to determine how to amortize the initial cost uniformly over 10 years; that is, we need to find the annual payments at 16% that are equivalent to the original cost. Using the annuity formula, we find that this corresponds to \$20,690 per year. Hence the annual worth of the project is  $\$25,000 - \$20,690 = \$4,310$ , which is positive; thus the investment is profitable. Note that if the purchase of the

If machine were financed at 16% over 10 years, the *actual* yearly net cash flows would correspond exactly to the annual worth.

### 3.3 Bond Details

Bonds represent by far the greatest monetary value of fixed-income securities and are, as a class, the most liquid of these securities. We devote special attention to bonds, both because of their practical importance as investment vehicles and because of their theoretical value, which will be exploited heavily in Chapter 4. We describe the general structure and trading mechanics of bonds in this section and then discuss in the following few sections some methods by which bonds are analyzed. Our description is intended to be an overview. Specific details are quite involved, and one must refer to specialized literature or to a brokerage firm for the exact features of any particular bond issue.

A **bond** is an obligation by the bond issuer to pay money to the bond holder according to rules specified at the time the bond is issued. Generally, a bond pays a specific amount, its **face value** or, equivalently, its **par value** at the date of maturity. Bonds generally have par values of even amounts, such as \$1,000 or \$ 10,000. In addition, most bonds pay periodic **coupon payments**. The term *coupon* is due to the fact that in the past actual coupons were attached to bond certificates. The bond holder would mail these to the agent of the issuer (usually a bank) one at a time, at specified dates, and the appropriate coupon payment would then be sent by return mail. These physical coupons are rare today, but the name remains. The last coupon date corresponds to the maturity date, so the last payment is equal to the face value plus the coupon value.

The coupon amount is described as a percentage of the face value. For example, a 9% coupon bond with a face value of \$1,000 will have a coupon of \$90 per year. However, the period between coupons may be less than a year. In the United States, coupon payments are generally made every 6 months, paying one-half of the coupon amount. This would be \$45 in our example.

The issuer of a bond initially sells the bonds to raise capital immediately, and then is obligated to make the prescribed payments. Usually bonds are issued with coupon rates close to the prevailing general rate of interest so that they will sell at close to their face value. However, as time passes, bonds frequently trade at prices different from their face values. While any two parties can agree on a price and execute a trade, the vast majority of bonds are sold either at auction (when originally issued) or through an exchange organization. The price is therefore determined by a market and thus may vary minute by minute.

**Treasury bills** are sold at a discount to their par value and make no coupon payment. They are sold at auction to major banks, which in turn sell them to other financial institutions. A sample of auction results for bills is shown in Table 3.3.

In the case of coupon bonds, the price quotations ignore **accrued interest**, which must be added to the price quoted in order to obtain the actual amount that must be paid for the bond. Suppose that a bond makes coupon payments every 6 months. If you purchase the bond midway through the coupon period, you will receive your first

**TABLE 3.3**  
**RECENT TREASURY BILL AUCTION RESULTS, APRIL 21, 2011**

Security term	Issue date	Maturity date	Discount rate	Investment rate	Price per \$100
4 weeks	04-21-2011	05-19-2011	0.030	0.030	99.997667
13 weeks	04-21-2011	07-21-2011	0.030	0.061	99.984833
26 weeks	04-21-2011	10-20-2011	0.110	0.112	99.944389
4 weeks	04-14-2011	05-12-2011	0.025	0.025	99.998056
13 weeks	04-14-2011	07-14-2011	0.050	0.051	99.987361
26 weeks	04-14-2011	10-13-2011	0.110	0.112	99.944389

*The discount rate is the yield based on the selling price and a 360-day year. The investment date is based on a 365- or 366-day year.*

coupon payment after only 3 months. You are getting extra interest—interest that was, in theory, earned by the previous owner. So you must pay the first 3 months' interest to the previous owner. This interest payment is made at the time of the sale, not when the next coupon payment is made, so this extra payment acts like an addition to the price. The accrued interest that must be paid to the previous owner is determined by a straight-line interpolation based on days. Specifically, the accrued interest (AI) is

$$AI = \frac{\text{number of days since last coupon}}{\text{number of days in current coupon period}} \times \text{coupon amount}.$$

**Example 3.5 (Accrued interest calculation)** Suppose we purchase on May 8 a U.S. Treasury bond that matures on August 15 in some distant year. The coupon rate is 9%. Coupon payments are made every February 15 and August 15. The accrued interest is computed by noting that there have been 83 days since the last coupon (in a leap year) and 99 days until the next coupon payment. Hence,

$$AI = \frac{83}{83 + 99} \times 4.50 = 2.05.$$

This 2.05 would be added to the quoted price, expressed as a percentage of the face value. For example, \$20.50 would be added to the bond if its face value were \$1,000.

## Quality Ratings

Although bonds offer a supposedly fixed-income stream, they are subject to default if the issuer has financial difficulties or falls into bankruptcy. To characterize the nature of this risk, bonds are rated by rating organizations. The two primary rating classifications are issued and published by Moody's and Standard & Poor's. Their classification schemes are shown in Table 3.4. U.S. Treasury securities are not rated, since they are considered to be essentially free of default risk.

**TABLE 3.4**  
**RATING CLASSIFICATIONS**

	<b>Moody's</b>	<b>Standard &amp; Poor's</b>
High grade	Aaa	AAA
	Aa	AA
Medium grade	A	A
	Baa	BBB
Speculative grade	Ba	BB
	B	B
Default danger	Caa	CCC
	Ca	CC
	C	C
		D

*Ratings reflect judgment of the likelihood that bond payments will be made as scheduled. Bonds with low ratings usually sell at lower prices than comparable bonds with high ratings. (Also see Section 17.3.)*

Bonds that are either high or medium grade are considered to be **investment grade**. Bonds that are in or below the speculative category are often termed **junk bonds**. Historically, the frequency of default has correlated well with the assigned ratings.

The assignment of a rating class by a rating organization is largely based on the issuer's financial status as measured by various financial ratios. For example, the ratio of debt to equity, the ratio of current assets to current liabilities, the ratio of cash flow to outstanding debt, as well as several others are used. The trend in these ratios is also considered important. (See Chapter 17 for more details.)

A bond with a low rating will have a lower price than a comparable bond with a high rating. Hence some people have argued that junk bonds may occasionally offer good value if the default risk can be diversified. A careful analysis of this approach requires explicit consideration of uncertainty, however.

## 3.4 Yield

A bond's yield is the interest rate implied by the payment structure. Specifically, it is the interest rate at which the present value of the stream of payments (consisting of the coupon payments and the final face-value redemption payment) is exactly equal to the current price. This value is termed more properly the **yield to maturity** (YTM) to distinguish it from other yield numbers that are sometimes used. Yields are always quoted on an annual basis.

It should be clear that the yield to maturity is just the internal rate of return of the bond at the current price. But when discussing bonds, the term *yield* is generally used instead.

Suppose that a bond with face value  $F$  makes  $m$  coupon payments of  $C/m$  each year and there are  $n$  periods remaining. The coupon payments sum to  $C$  within a year. Suppose also that the current price of the bond is  $P$ . Then the yield to maturity is the value of  $\lambda$  such that

$$P = \frac{F}{[1 + (\lambda/m)]^n} + \sum_{k=1}^n \frac{C/m}{[1 + (\lambda/m)]^k}. \quad (3.1)$$

This value of  $\lambda$ , the yield to maturity, is the interest rate implied by the bond when interest is compounded  $m$  times per year. Note that the first term in equation (3.1) is the present value of the face-value payment. The  $k$ th term in the summation is the present value of the  $k$ th coupon payment  $C/m$ . The sum of the present values, based on a nominal interest rate of  $\lambda$ , is set equal to the bond's price.

The summation in (3.1) can be collapsed by use of the general value formula for annuities in the previous section, since this sum represents the present value of the equal coupon payments of  $C/m$ . The collapsed form is highlighted here:

**Bond price formula** *The price of a bond, having exactly  $n$  coupon periods remaining to maturity and a yield to maturity of  $\lambda$ , satisfies*

$$P = \frac{F}{[1 + (\lambda/m)]^n} + \frac{C}{\lambda} \left\{ 1 - \frac{1}{[1 + (\lambda/m)]^n} \right\}, \quad (3.2)$$

where  $F$  is the face value of the bond,  $C$  is the yearly coupon payment, and  $m$  is the number of coupon payments per year.

Equation (3.2) must be solved for  $\lambda$  to determine the yield. This cannot be done by hand except for very simple cases. It should be clear that the terms in equation (3.2) are the familiar terms giving the present value of a single future payment and of an annuity. However, to determine  $\lambda$  one must do more than just evaluate these expressions. One must adjust  $\lambda$  so that equation (3.2) is satisfied. As in any calculation of internal rate of return, this generally requires an iterative procedure, easily carried out by a computer. There are, however, specialized calculators and bond tables devised for this purpose, which are used by bond dealers and other professionals. Spreadsheet packages also typically have built-in bond formulas.

The formulas discussed here assume that there is an exact number of coupon periods remaining to the maturity date. The price–yield formula requires adjustment for dates between coupon payment dates.

## Qualitative Nature of Price–Yield Curves

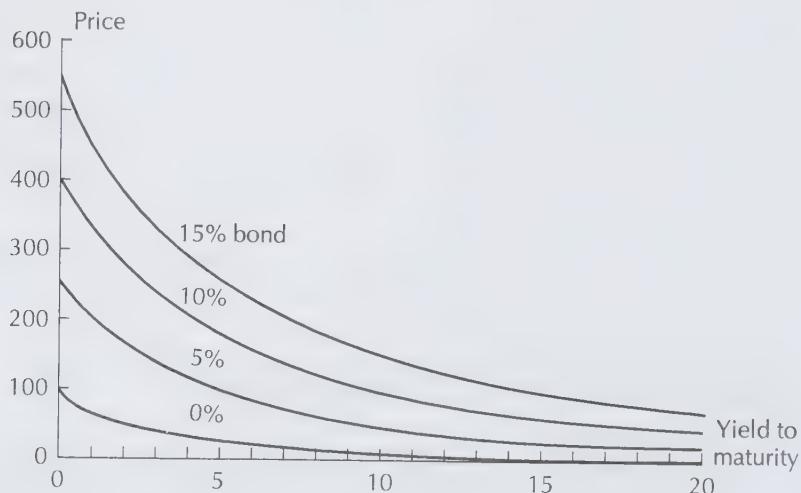
Although the bond equation is complex, it is easy to obtain a qualitative understanding of the relationship between price, yield, coupon, and time to maturity. This qualitative understanding helps motivate the ideas underlying bond portfolio construction and, specifically, leads to an understanding of the interest rate risk properties of bonds. The following examples should be studied with an eye toward obtaining this kind of understanding.

As a general rule, the yields of various bonds track one another and the prevailing interest rates of other fixed-income securities quite closely. After all, most people would not buy a bond with a yield of 6% when bank CDs are offering 10%. The general interest rate environment exerts a force on every bond, urging its yield to conform to that of other bonds. However, the only way that the yield of a bond can change is for the bond's price to change. So as yields move, prices move correspondingly. But the price change required to match a yield change varies with the structure of the bond (its coupon rate and its maturity). So as the yields of various bonds move more or less in harmony, their prices move by different amounts. To understand bonds, it is important to understand this relation between the price and the yield. For a given bond, this relationship is shown pictorially by the **price–yield curve**.

Examples of price–yield curves are shown in Figure 3.3. Here the price, as a percentage of par, is shown as a function of YTM expressed in percentage terms. Let us focus on the bond labeled 10%. This bond has a 10% coupon, which means 10% of the face value is paid each year (or 5% every 6 months), and it has 30 years to maturity. The price–yield curve shows how yield and price are related.

The first obvious feature of the curve is that it has negative slope; that is, price and yield have an inverse relation. If yield goes up, price goes down. If I am to obtain a higher yield on a fixed stream of received payments, the price I pay for this stream must be lower. This is a fundamental feature of bond markets. When people say “the bond market went down,” they mean that interest rates went up.

Some points on the curve can be calculated by inspection. First, suppose that  $YTM = 0$ . This means that the bond is priced as if it offered no interest. Within the framework of this bond, money in the future is not discounted. In that case, the present value of the bond is just equal to the sum of all payments: here coupon payments of

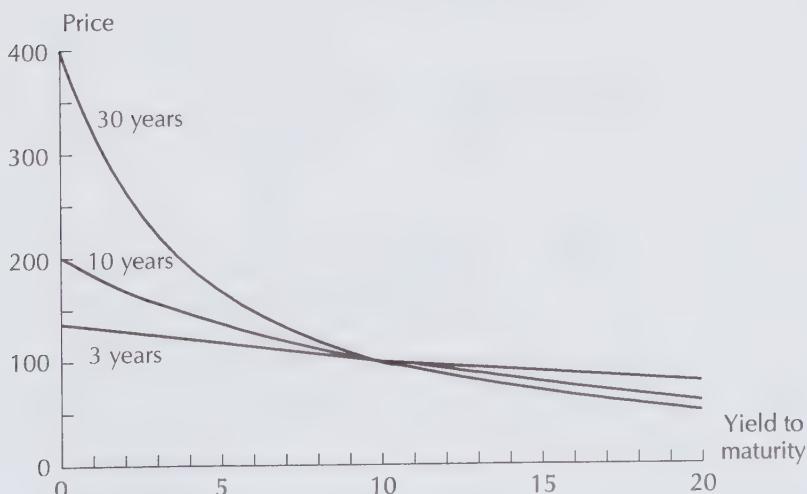


**FIGURE 3.3 Price–yield curves and coupon rate.** All bonds shown have a maturity of 30 years and the coupon rates indicated on the respective curves. Prices are expressed as a percentage of par.

10 points each year for 30 years, giving 300, plus the 100% of par value received at maturity, for a total of 400. This is the value of the bond at zero yield. Second, suppose that  $YTM = 10\%$ . Then the value of the bond is equal to the par value. The reason for this is that each year the coupon payment just equals the 10% yield expected on the investment. The value remains at 100 every year. The bond is like a loan where the interest on the principal is paid each year and hence the principal remains constant. In this situation, where the yield is exactly equal to the coupon rate, the bond is termed a **par bond**. In addition to these two specific points on the price–yield curve, we can deduce that the price of the bond must tend toward zero as the yield increases—large yields imply heavy discounting, so even the nearest coupon payment has little present value. Overall, the shape of the curve is **convex** since it bends toward the origin and out toward the horizontal axis. Just given the two points and this rough knowledge of shape, it is possible to sketch a reasonable approximation to the true curve.

Let us briefly examine another one of the curves, say, the 15% bond. The price at  $YTM = 0$  is  $15 \times 30 + 100 = 550$ , and the par point of 100 is at 15%. We see that with a fixed maturity date, the price–yield curve rises as the coupon rate increases.

Now let us consider the influence of the time to maturity. Figure 3.4 shows the price–yield curves for three different bonds. Each of these bonds has a 10% coupon rate, but they have different maturities: 30 years, 10 years, and 3 years. All of these bonds are at par when the yield is 10%; hence the three curves all pass through the common par point. However, the curves pivot upward around that point by various amounts, depending on the maturity. The values at  $YTM = 0$  can be found easily, as before, by simply summing the total payments. The main feature is that as the maturity is increased, the price–yield curve becomes steeper, essentially pivoting about the par



**FIGURE 3.4 Price–yield curves and maturity.** The price–yield curve is shown for three maturities. All bonds have a 10% coupon.

**TABLE 3.5**  
**PRICES OF 9% COUPON BONDS**

<b>Time to maturity</b>	<b>Yield</b>				
	<b>5%</b>	<b>8%</b>	<b>9%</b>	<b>10%</b>	<b>15%</b>
1 year	103.85	100.94	100.00	99.07	94.61
5 years	117.50	104.06	100.00	96.14	79.41
10 years	131.18	106.80	100.00	93.77	69.42
20 years	150.21	109.90	100.00	91.42	62.22
30 years	161.82	111.31	100.00	90.54	60.52

*The prices of long-maturity bonds are more sensitive to yield changes than are the prices of bonds of short maturity.*

point. This increased steepness is an indication that longer maturities imply greater sensitivity of price to yield.

The price–yield curve is important because it describes the interest rate risk associated with a bond. For example, suppose that you purchased the 10% bond illustrated in Figure 3.3 at par (when the yield was 10%). It is likely that all bonds of maturity approximately equal to 30 years would have yields of 10%, even though some might not be at par. Then 10% would represent the market rate for such bonds. Now suppose that market conditions change and the yield on your bond increases to 11%. The price of your bond will drop to 91.28. This represents an 8.72% change in the value of your bond. It is good to consider the possibility of such a change when purchasing this bond. For example, with a 3-year 10% par bond, if the yield rose to 11%, the price would drop only to 97.50, and hence the interest rate risk is lower with this bond. Of course if yields *decreased*, you would *profit* by similar amounts.

Bond holders are subject to yield risk in the sense described: if yields change, bond prices also change. This is an immediate risk, affecting the near-term value of the bond. You may, of course, continue to hold the bond and thereby continue to receive the promised coupon payments and the face value at maturity. This cash flow stream is not affected by interest rates. (That is after all why the bond is classified as a fixed-income security.) But if you plan to sell the bond before maturity, the price will be governed by the price–yield curve.

Table 3.5 displays the price–yield relation in tabular form for bonds with a 9% coupon rate. It is easy to see that the bond with 30-year maturity is much more sensitive to yield changes than the bond with 1-year maturity.

It is the quantification of this risk that underlies the importance of the price–yield relation. Our rough qualitative understanding is important. The next sections develop additional tools for studying this risk.

## Other Yield Measures

Other measures of yield, aside from yield to maturity, are used to gain additional insight into a bond's properties. For example, one important yield measure is **current**

**yield (CY)**, which is defined as

$$\text{CY} = \frac{\text{annual coupon payment}}{\text{bond price}} \times 100.$$

The current yield gives a measure of the annual return of the bond. For instance, consider a 10%, 30-year bond. If it is selling at par (that is, at 100), then the current yield is 10, which is identical to the coupon rate and to the yield to maturity. If the same bond were selling for 90, then  $\text{CY} = 10/90 = 11.11$  while  $\text{YTM} = 11.16$ .

Another measure, used if the bond is callable after some number of years, is the **yield to call (YTC)**, which is defined as the internal rate of return calculated assuming that the bond is in fact called at the earliest possible date.

There are several other yield measures that account for sinking funds, principal payments, and other features.

## 3.5 Duration

Everything else being equal, bonds with long maturities have steeper price–yield curves than bonds with short maturities. Hence the prices of **long bonds** are more sensitive to interest rate changes than those of **short bonds**. This is shown clearly in Table 3.5. However, this is only a rough rule of thumb. Maturity itself does not give a complete quantitative measure of interest rate sensitivity.

Another measure of time length termed **duration** does give a direct measure of interest rate sensitivity. This section describes this measure.

The duration of a fixed-income instrument is a weighted average of the times that payments (cash flows) are made. The weighting coefficients are the present values of the individual cash flows.

We can write out this definition more explicitly. Suppose that cash flows are received at times  $t_0, t_1, t_2, \dots, t_n$ . Then the duration of this stream is

$$D = \frac{\text{PV}(t_0)t_0 + \text{PV}(t_1)t_1 + \text{PV}(t_2)t_2 + \cdots + \text{PV}(t_n)t_n}{\text{PV}}$$

In this formula the expression  $\text{PV}(t_k)$  denotes the present value of the cash flow that occurs at time  $t_k$ . The term PV in the denominator is the total present value, which is the sum of the individual  $\text{PV}(t_k)$  values.

The expression for  $D$  is indeed a weighted average of the cash flow times. Hence  $D$  itself has units of time. When the cash flows are all nonnegative, as they are for a bond already owned (so that the purchase is not included in the cash flow), then it is clear that  $t_0 \leq D \leq t_n$ . Duration is a time intermediate between the first and last cash flows.

Clearly, a zero-coupon bond, which makes only a final payment at maturity, has a duration equal to its maturity date. Nonzero-coupon bonds have durations strictly less than their maturity dates. This shows that duration can be viewed as a generalized maturity measure. It is an average of the maturities of all the individual payments.

## Interest Duration

The preceding definition does not state explicitly how the present value is to be calculated. One natural choice is to use the prevailing interest rate defined for continuous-time analysis; that is, the present value of a cash flow  $x(t)$  at time  $t$  is  $\text{PV}(t) = e^{-rt}x(t)$ , where  $r$  is the interest rate. The total present value of a cash flow stream is then

$$\text{PV} = e^{-rt_0}x(t_0) + e^{-rt_1}x(t_1) + \cdots + e^{-rt_n}x(t_n).$$

Differentiating with respect to  $r$  produces

$$\begin{aligned}\frac{d\text{PV}}{dr} &= e^{-rt_0}(t_0)t_0 + e^{-rt_1}x(t_1)t_1 + \cdots + e^{-rt_n}x(t_n)t_n \\ &= -D\text{PV}.\end{aligned}$$

In other words,

$$\frac{1}{\text{PV}} \frac{d\text{PV}}{dr} = -D.$$

This shows how this type of duration measures directly the relative sensitivity of a present value with respect to changes in the interest rate.

## Macaulay Duration

In the study of bonds, it is common to use yield rather than an interest rate. After all, it is yield that relates directly to price. In this approach the general duration formula becomes the Macaulay duration.

Specifically, suppose a financial instrument makes payments  $m$  times per year, with the payment in period  $k$  being  $c_k$ , and there are  $n$  periods remaining. The **Macaulay duration**  $D$  is defined as

$$D = \frac{\sum_{k=1}^n (k/m)c_k/[1 + (\lambda/m)]^k}{\text{PV}}$$

where  $\lambda$  is the yield to maturity and

$$\text{PV} = \sum_{k=1}^n \frac{c_k}{[1 + (\lambda/m)]^k}.$$

Note that the factor  $k/m$  in the numerator of the formula for  $D$  is time, measured in years. In this chapter we always use the Macaulay duration (or a slight modification of it), and hence we do not give it a special symbol, but denote it by  $D$ , the same as in the general definition of duration.

**Example 3.6 (A short bond)** Consider a 7% bond with 3 years to maturity. Assume that the bond is selling at 8% yield. We can find the value and the Macaulay duration by the simple spreadsheet layout shown in Figure 3.5. The duration is 2.753 years.

**FIGURE 3.5 Layout for calculating duration.** Present values of payments are calculated in column D. Dividing these by the total present value gives the weights shown in column E. The duration is obtained using this weighted average of the payment times.

A	B	C	D	E	F
Year	Payment	Discount factor (@ 8%)	Present value of payment ( $B \times C$ )	Weight (D/Price)	$A \times E$
.5	3.5	.962	3.365	.035	.017
1	3.5	.925	3.236	.033	.033
1.5	3.5	.889	3.111	.032	.048
2	3.5	.855	2.992	.031	.061
2.5	3.5	.822	2.877	.030	.074
3	103.5	.790	81.798	.840	2.520
Sum			97.379	1.000	2.753
			Price		Duration

## Explicit Formula\*

In the case where all coupon payments are identical (which is the normal case for bonds) there is an explicit formula for the sum of the series that appears in the numerator of the expression for the Macaulay duration. We skip the algebra here and just give the result.

**Macaulay duration formula** *The Macaulay duration for a bond with a coupon rate  $c$  per period, yield  $y$  per period,  $m$  periods per year, and exactly  $n$  periods remaining, is*

$$D = \frac{1+y}{my} - \frac{1+y+n(c-y)}{mc[(1+y)^n - 1] + my}. \quad (3.3)$$

**Example 3.7 (Duration of a 30-year par bond)** Consider a 10%, 30-year bond with 6-month coupons. Let us assume that it is at par; that is, the yield is 10%. At par,  $c = y$ , and equation (3.3) reduces to

$$D = \frac{1+y}{my} \left[ 1 - \frac{1}{(1+y)^n} \right]$$

Hence,

$$D = \frac{1.05}{.1} \left[ 1 - \frac{1}{(1.05)^{60}} \right] = 9.938.$$

## Qualitative Properties of Duration\*

The duration of a coupon-paying bond is always less than its maturity, but often it is surprisingly short. An appreciation for the relation between a bond's duration and other parameters of the bond can be obtained by examination of Table 3.6. In

**TABLE 3.6**  
**DURATION OF A BOND YIELDING 5% AS FUNCTION OF Maturity AND COUPON RATE**

Years to maturity	Coupon rate			
	1%	2%	5%	10%
1	.997	.995	.988	.977
2	1.984	1.969	1.928	1.868
5	4.875	4.763	4.485	4.156
10	9.416	8.950	7.989	7.107
25	20.164	17.715	14.536	12.754
50	26.666	22.284	18.765	17.384
100	22.572	21.200	20.363	20.067
$\infty$	20.500	20.500	20.500	20.500

*Duration does not increase appreciably with maturity. In fact, with fixed yield, duration increases only to a finite limit as maturity is increased.*

this table the yield is held fixed at 5%, but various maturities and coupon rates are considered. This procedure approximates the situation of looking through a list of available bonds at a time when all yields hover near 5%. Within a given class (say, government securities) the available bonds then differ mainly by these two parameters.

One striking feature of this table is that as the time to maturity increases to infinity, the durations do *not* also increase to infinity, but instead tend to a finite limit that is independent of the coupon rate. (See Exercise 18.) Another feature of the table is that the durations do not vary rapidly with respect to the coupon rate. The fact that the yield is held constant tends to cancel out the influence of the coupons.

A general conclusion is that very long durations (of, say, 20 years or more) are achieved only by bonds that have both very long maturities and very low coupon rates.

## Duration and Sensitivity

Duration is useful because it measures directly the sensitivity of price to changes in yield. This follows from a simple expression for the derivative of the present value expression.

In the case where payments are made  $m$  times per year and yield is based on those same periods, we have

$$PV_k = \frac{c_k}{[1 + (\lambda/m)]^k}.$$

The derivative with respect to  $\lambda$  is

$$\frac{dPV_k}{d\lambda} = \frac{-(k/m)c_k}{[1 + (\lambda/m)]^{k+1}} = -\frac{k/m}{1 + (\lambda/m)}PV_k.$$

We now apply this to the expression for price,

$$P = \sum_{k=1}^n PV_k.$$

Here we have used the fact that the price is equal to the total present value at the yield (by definition of yield). We find that

$$\frac{dP}{d\lambda} = \sum_{k=1}^n \frac{dPV_k}{d\lambda} = - \sum_{k=1}^n \frac{(k/m)PV_k}{1 + (\lambda/m)} = - \frac{1}{1 + (\lambda/m)} DP \equiv -D_M P. \quad (3.4)$$

The value  $D_M$  is called the **modified duration**. It is the usual duration modified by the extra term in the denominator. Note that  $D_M \approx D$  for large values of  $m$  or small values of  $\lambda$ . We highlight this important sensitivity relation:

**Price sensitivity formula** *The derivative of price  $P$  with respect to yield  $\lambda$  of a fixed-income security is*

$$\frac{dP}{d\lambda} = -D_M P \quad (3.5)$$

where  $D_M = D/[1 + (\lambda/m)]$  is the modified duration.

It is perhaps most revealing to write equation (3.5) as

$$\frac{1}{P} \frac{dP}{d\lambda} = -D_M.$$

The left side is then the relative change in price (or the fractional change). Hence  $D_M$  measures the relative change in a bond's price directly as  $\lambda$  changes.

By using the approximation  $dP/d\lambda \approx \Delta P/\Delta\lambda$ , equation (3.5) can be used to estimate the change in price due to a small change in yield (or vice versa). Specifically, we would write

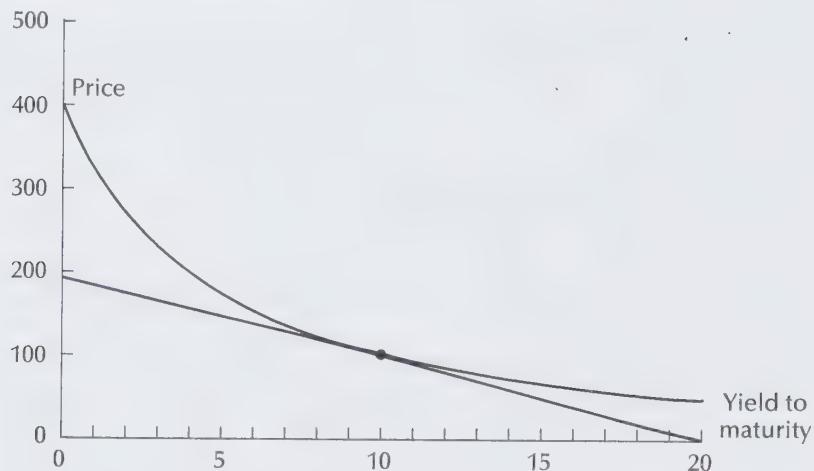
$$\Delta P \approx -D_M P \Delta\lambda.$$

This gives explicit values for the impact of yield variations.

**Example 3.8 (A 10% bond)** The price–yield curve for a 30-year, 10% coupon bond is shown in Figure 3.6. As computed earlier, the duration of this bond at the par point (where price = 100) is  $D = 9.94$ . Hence  $D_M = 9.94/1.05 = 9.47$ . The slope of the price–yield curve at that point is, according to equation (3.5), equal to  $dP/d\lambda = -947$ . A line with this slope can be placed tangent to the price–yield curve at the point where the duration was calculated, as shown in Figure 3.6. This line provides a straight-line approximation to the curve for nearby points. For example, if the yield changes to 11%, we can estimate the change in price as

$$\Delta P = -D_M 100 \Delta\lambda = -947 \times .01 = -9.47.$$

Hence  $P \approx 90.53$ .



**FIGURE 3.6 Price–yield curve and slope.** The slope of the line tangent to the curve at  $P$  is  $-D_M P$ .

**Example 3.9 (A zero-coupon bond)** Consider a 30-year zero-coupon bond. Suppose its current yield is 10%. Then we have  $D = 30$  and  $D_M \approx 27$ . Suppose that yields increase to 11%. According to equation (3.5), the relative price change is approximately equal to 27%. This is a very large loss in value. Because of their long durations, zero-coupon bonds have very high interest rate risk.

## Duration of a Portfolio

Suppose that a portfolio of several bonds of different maturities is assembled. This portfolio acts like a master fixed-income security: it receives periodic payments, but due to the different maturities, the payments may not be of equal magnitude. What can we say about the duration of this portfolio?

First, suppose that all the bonds have the same yield. (This is usually approximately true, since yields tend to track each other closely, if not exactly.) The duration of the portfolio is then just a weighted sum of the durations of the individual bonds—with weighting coefficients proportional to individual bond prices. We can easily verify this for a portfolio that is the sum of two bonds A and B. The durations are

$$D^A = \frac{\sum_{k=0}^n t_k PV_k^A}{P^A}$$

$$D^B = \frac{\sum_{k=0}^n t_k PV_k^B}{P^B}.$$

Hence,

$$P^A D^A + P^B D^B = \sum_{k=0}^n t_k (\text{PV}_k^A + \text{PV}_k^B)$$

which gives, upon division by  $P = P^A + P^B$ ,

$$D = \frac{P^A D^A}{P} + \frac{P^B D^B}{P}$$

as the duration of the portfolio. Therefore  $D$  is a weighted average of the durations of the individual bonds, with the weight of a bond's duration being proportional to that bond's price. The result easily extends to a portfolio containing several bonds.

**Duration of a portfolio** Suppose there are  $m$  fixed-income securities with prices and durations of  $P_i$  and  $D_i$ , respectively,  $i = 1, 2, \dots, m$ , all computed at a common yield. The portfolio consisting of the aggregate of these securities has price  $P$  and duration  $D$ , given by

$$P = P_1 + P_2 + \cdots + P_m$$

$$D = w_1 D_1 + w_2 D_2 + \cdots + w_m D_m,$$

where  $w_i = P_i/P$ ,  $i = 1, 2, \dots, m$ .

The duration of a portfolio measures the interest rate sensitivity of that portfolio, just as normal duration measures it for a single bond. That is, if the yield changes by a small amount, the total value of the portfolio will change approximately by the amount predicted by the equation relating prices to (modified) duration.

If the bonds composing a portfolio have different yields, the composite duration as defined can still be used as an approximation. In this case a single yield must be chosen—perhaps the average. Then present values can be calculated with respect to this single yield value, although these present values will not be exactly equal to the prices of the bonds. The weighted average duration, calculated as shown, will give the sensitivity of the overall present value to a change in the yield figure that is used.

## 3.6 Immunization

We now have the concepts and tools necessary to solve a problem of major practical value, namely, the structuring of a bond portfolio to protect against interest rate risk. This procedure is termed **immunization** because it “immunizes” the portfolio value against interest rate changes. The procedure, as well as its refinements, is in fact one of the most (if not *the* most) widely used analytical techniques of investment science, shaping portfolios consisting of billions of dollars of fixed-income securities held by pension funds, insurance companies, and other financial institutions.

Before describing the procedure, let us more fully consider its purpose. A portfolio cannot be structured meaningfully without a statement of its purpose. The purpose helps define the character of risk one is willing to assume. For example, suppose that you wish to invest money now that will be used next year for a major household

expense. If you invest in 1-year Treasury bills, you know exactly how much money these bills will be worth in a year, and hence there is little risk relative to your purpose. If, on the other hand, you invested in a 10-year zero-coupon bond, the value of this bond a year from now would be quite variable, depending on what happens to interest rates during the year. This investment has high risk relative to your purpose. The situation would be reversed if you were saving the money to pay off an obligation that was due in 10 years. Then the 10-year zero-coupon bond would provide completely predictable results, but the 1-year Treasury bill would impose **reinvestment risk** since the proceeds would have to be reinvested after 1 year at the then prevailing rate (which could be considerably lower than the current rate).

Suppose now that you face a series of cash obligations and you wish to acquire a portfolio that you will use to pay these obligations as they arise. (This is the sort of problem faced by life insurance companies.) One way to do this is to purchase a set of zero-coupon bonds that have maturities and face values exactly matching the separate obligations. However, this simple technique may not be feasible if corporate bonds are used, since there are few corporate zero-coupon bonds. (You may wish to use corporate bonds because they offer higher yields than U.S. Treasury bonds.) If perfect matching is not possible, you may instead acquire a portfolio having a value equal to the present value of the stream of obligations. You can sell some of your portfolio whenever cash is needed to meet a particular obligation; or if your portfolio delivers more cash than needed at a given time (from coupon or face value payments), you can buy more bonds. If the yield does not change, the value of your portfolio will, throughout this process, continue to match the present value of the remaining obligations. Hence you will meet the obligations exactly.

A problem with this present-value-matching technique arises if the yields change. The value of your portfolio and the present value of the obligation stream will both change in response, but probably by amounts that differ from one another. Your portfolio will no longer be matched.

Immunization solves this problem—at least approximately—by matching durations as well as present values. If the duration of the portfolio matches that of the obligation stream, then the cash value of the portfolio and the present value of the obligation stream will respond identically (to first order) to a change in yield. Specifically, if yields increase, the present value of the asset portfolio will decrease, but the present value of the obligation will decrease by approximately the same amount; so the value of the portfolio will still be adequate to cover the obligation. The process is best explained through an example.

**Example 3.10 (The X Corporation)** The X Corporation has an obligation to pay \$1 million in 10 years. It wishes to invest money now that will be sufficient to meet this obligation.

The purchase of a single zero-coupon bond would provide one solution; but such zeros are not always available in the required maturities. We assume that none are available for this example. Instead the X Corporation is planning to select from the three corporate bonds shown in Table 3.7.

These bonds all have the same yield of 9%, and this rate is used in all calculations. The X Corporation first considers using bonds 2 and 3 to construct its portfolio.

**TABLE 3.7**  
**BOND CHOICES**

	Rate	Maturity	Price	Yield
Bond 1	6%	30 yr	69.04	9.00%
Bond 2	11%	10 yr	113.01	9.00%
Bond 3	9%	20 yr	100.00	9.00%

*Three bonds are considered for the X Corporation's immunized portfolio.*

**TABLE 3.8**  
**IMMUNIZATION RESULTS**

	Percent yield		
	9.0	8.0	10.0
Bond 1			
Price	69.04	77.38	62.14
Shares	4,241.00	4,241.00	4,241.00
Value	292,798.64	328,168.58	263,535.74
Bond 2			
Price	113.01	120.39	106.23
Shares	1,078.00	1,078.00	1,078.00
Value	121,824.78	129,780.42	114,515.94
Obligation			
Value	414,642.86	456,386.95	376,889.48
Surplus	−19.44	1,562.05	1,162.20

*The net surplus of portfolio value minus obligation value remains approximately equal to zero even if yields change.*

As a first step it calculates the durations and finds  $D_2 = 6.54$  and  $D_3 = 9.61$ , respectively. This is a serious problem! The duration of the obligation is obviously 10 years, and there is no way to attain that with a weighted average of  $D_2$  and  $D_3$  using positive weights. A bond with a longer duration is required. Therefore the X Corporation decides to use bonds 1 and 2. It is found that  $D_1 = 11.44$ . (Note that, consistent with the discussion on the qualitative nature of durations, it is quite difficult to obtain a long duration when the yield is 9%—a long maturity and a low coupon are required.) Fortunately  $D_1 > 10$ , and hence bonds 1 and 2 will work.

Next the present value of the obligation is computed at 9% interest. This is  $PV = \$414,643$ . The immunized portfolio is found by solving the two equations

$$V_1 + V_2 = PV$$

$$D_1 V_1 + D_2 V_2 = 10PV$$

for the amounts of money  $V_1$  and  $V_2$  to be invested in the two bonds. The first equation states that the total value of the portfolio must equal the total present value of the obligation. The second states that the duration of the portfolio must equal the duration

(10 years) of the obligation. (The modification terms in the durations can be omitted, since they cancel each other out.) The solution to these equations is  $\dot{V}_1 = \$292,788.73$  and  $V_2 = \$121,854.27$ . The number of bonds to be purchased is then found by dividing each value by the respective bond price. (We assume a face value of \$100.) These numbers are then rounded to integers to define the portfolio.

The results are shown in Table 3.8. Note that, except for rounding error, the present value of the portfolio does indeed equal that of the obligation. Furthermore, at different yields (8% and 10% are shown) the value of the portfolio is still approximately equal to that of the obligation. In fact, due to the structure of the price–yield curve, the portfolio value will always exceed the value of the obligation in both cases. (See Exercise 20.)

Immunization provides protection against changes in yield. If the yield changes immediately after purchase of the portfolio, the new value of the portfolio will, in theory, still approximately match the new value of the future obligation. However, once the yield does change, the new portfolio will not be immunized at the new rate. It is therefore desirable to **rebalance**, or reimmunize, the portfolio from time to time. Also, in practice more than two bonds would be used, partly to diversify default risk if the bonds included are not U.S. Treasury bonds.

Immunization is a clever idea, but it suffers some shortcomings, at least in this simple form. The method assumes that all yields are equal, whereas in fact they usually are not. Indeed it is quite unrealistic to assume that both long- and short-duration bonds can be found with identical yields. Usually long bonds have somewhat higher yields than short bonds. Furthermore, when yields change, it is unlikely that the yields on all bonds will change by the same amount; hence rebalancing would be difficult. We shall consider some important extensions of immunization in the next chapter, and in Chapter 5 we shall consider other approaches to bond portfolio construction. Overall, however, the technique given here is surprisingly practical.

## 3.7 Convexity\*

Modified duration measures the relative slope of the price–yield curve at a given point. As we have seen, this leads to a straight-line approximation to the price–yield curve that is useful both as a means of assessing risk and as a procedure for controlling it.

An even better approximation can be obtained by including a second-order (or quadratic) term. This second-order term is based on **convexity**, which is the relative curvature at a given point on the price curve.

For the case where the continuous-time interest rate  $r$  is used to measure present value, convexity is defined as

$$C = \frac{1}{PV} \frac{d^2PV}{dr^2}.$$

This leads to the simple formula

$$C = \sum_{k=1}^n \frac{1}{PV} [e^{-rt_k} x(t_k) t_k^2].$$

In other words, in this case convexity is a weighted average of the squares of the cash flow times, the weights being the present values.

In the study of bonds, it is more common to use yield instead of interest rate, since by definition yield is exactly related to price. In this case, convexity is the value of  $C$ , defined as

$$C = \frac{1}{P} \frac{d^2 P}{d\lambda^2},$$

which can be expressed in terms of the cash flow stream as

$$C = \frac{1}{P} \sum_{k=1}^n \frac{d^2 PV_k}{d\lambda^2}.$$

Assuming  $m$  coupons (and  $m$  compounding periods) per year, we have

$$C = \frac{1}{P[1 + (\lambda/m)]^2} \sum_{k=1}^n \frac{k(k+1)}{m^2} \frac{c_k}{[1 + (\lambda/m)]^k}.$$

Note that convexity has units of time squared. Convexity is the weighted average of  $t_k t_{k+1}$  where, like for duration, the weights are proportional to the present values of the corresponding cash flows. Then the result is modified by the factor  $1/[1 + (\lambda/m)]^2$ . An explicit formula can be derived for the case of equal-valued coupon payments.

Suppose that at a price  $P$  and a corresponding yield  $\lambda$ , the modified duration  $D_M$  and the convexity  $C$  are calculated. Then if  $\Delta\lambda$  is a small change in  $\lambda$  and  $\Delta P$  is the corresponding change in  $P$ , we have

$$\Delta P \approx -D_M P \Delta\lambda + \frac{PC}{2} (\Delta\lambda)^2.$$

This is the second-order approximation to the price–yield curve. Convexity can be used to improve immunization in the sense that, compared to ordinary immunization, a closer match of asset portfolio value and obligation value is maintained as yields vary. To account for convexity in immunization, one structures a portfolio of bonds such that its present value, its duration, and its convexity match those of the obligation. Generally, at least three bonds are required for this purpose.

## 3.8 Summary

Fixed-income securities are fundamental investment instruments, which are part of essentially every investment portfolio, and which reflect the market conditions for interest rates directly.

There are numerous kinds of fixed-income securities, designed for various investment and business purposes. However, the vast bulk of money in fixed-income securities is committed to mortgages and bonds.

Many fixed-income securities make periodic payments to the owner of the security. This is true, in particular, for mortgages, loans, annuities, and bonds. In

the case of bonds, these payments are usually made every 6 months and are termed coupon payments.

Usually the periodic payments associated with a fixed-income security are of equal magnitude, and there is an important formula relating the payment amount  $A$ , the principal value of the security  $P$ , the single-period interest rate  $r$ , and the number of payment periods  $n$ :

$$P = \frac{A}{r} \left[ 1 - \frac{1}{(1+r)^n} \right].$$

This single formula can be used to evaluate most annuities, mortgages, and bonds, and it can be used to amortize capital expenses over time.

Bonds are the most important type of fixed-income security for general investment purposes. Important reference bonds are U.S. Treasury securities—bills, notes, and bonds—of various maturities and coupon values. These bonds are considered to be default free and thus carry prices that are somewhat higher than corporate securities with similar coupon rates and maturities.

There are many variations to the generic coupon bond—call features, sinking fund bonds, bonds whose coupon rates are tied to economic indices, and so forth. In addition, municipal bonds receive special tax treatment.

A special feature of bonds is that the buyer must usually pay accrued interest in addition to the quoted price. This accrued interest is compensation to the previous owner for the coupon interest that has been earned since the last coupon payment.

Bonds are frequently analyzed by computing the yield to maturity. This is the annual interest rate that is implied by the current price. It is the interest rate that makes the present value of the promised bond payments equal to the current bond price. This calculation of yield can be turned around: the price of a bond can be found as a function of the yield. This is the price–yield relation which, when plotted, produces the price–yield curve.

The slope of the price–yield curve is a measure of the sensitivity of the price to changes in yield. Since yields tend to track the prevailing interest rate, the slope of the price–yield curve is therefore a measure of the interest rate risk associated with a particular bond. As a general rule, long bonds have greater slope than short bonds, and thus long bonds have greater interest rate risk. A normalized version of the slope—the slope divided by the current bond price—is given by the (modified) duration of the bond. Hence duration (or, more exactly, modified duration) is a convenient measure of interest rate risk.

Immunization is the process of constructing a portfolio that has, to first order, no interest rate risk. The process is frequently applied by institutions, such as insurance companies and pension funds, that have large future payment obligations. They wish to prepare for these obligations by making appropriate investments in fixed-income securities. A portfolio is immunized if its present value is equal to that of the stream of obligations and if its duration matches that of the obligation. In other words, the net portfolio, consisting of the obligation stream and the fixed-income assets, has zero present value and zero duration.

## Exercises

- (Amortization) A debt of \$25,000 is to be amortized over 7 years at 7% interest. What value of monthly payments will achieve this?
- (Cycles and annual worth  $\diamond$ ) Given a cash flow stream  $X = (x_0, x_1, x_2, \dots, x_n)$ , a new stream  $X_\infty$  of infinite length is made by successively repeating the corresponding finite stream. The interest rate is  $r$ . Let  $P$  and  $A$  be the present value and the annual worth, respectively, of stream  $X$ . Finally, let  $P_\infty$  be the present value of stream  $X_\infty$ . Find  $A$  in terms of  $P_\infty$  and conclude that  $A$  can be used as well as  $P_\infty$  for evaluation purposes.
- (Uncertain annuity  $\diamond$ ) Gavin's grandfather, Mr. Jones, has just turned 90 years old and is applying for a lifetime annuity that will pay \$10,000 per year, starting 1 year from now, until he dies. He asks Gavin to analyze it for him. Gavin finds that according to statistical summaries, the chance (probability) that Mr. Jones will die at any particular age is as follows:

age	90	91	92	93	94	95	96	97	98	99	100	101
probability	.07	.08	.09	.10	.10	.10	.10	.10	.10	.07	.05	.04

Then Gavin (and you) answer the following questions:

- What is the life expectancy of Mr. Jones?
  - What is the present value of an annuity at 8% interest that has a lifetime equal to Mr. Jones's life expectancy? (For an annuity of a nonintegral number of years, use an averaging method.)
  - What is the expected present value of the annuity?
- (APR) For the mortgage listed second in Table 3.1 what are the total fees?
  - (Mortgage restructuring) An investor purchased a small apartment building for \$250,000. She made a down payment of \$50,000 and financed the balance with a 30-year, fixed-rate mortgage at 12% annual interest, compounded monthly. For exactly 20 years she has made equal-sized monthly payments as required by the terms of the loan. Now she has the opportunity to restructure the mortgage by refinancing the balance. She could borrow the current balance, pay off the original loan, and assume a new loan for the balance. (No points or any other charges are involved in the transaction.) The new loan is a 20-year, fixed-rate loan at 9%, compounded monthly, to be paid in equal monthly installments. Suppose she has a risk-free savings account that pays 5%, compounded monthly. Should she restructure the mortgage?
  - (Simple cash flow) Assume we are at period 0. The current interest is  $r$ . Define  $x = \frac{1}{1+r}$ .

- (a) Derive the present value  $S_n$  of the following cash flow in terms of  $x$  and  $n$ :

Period	0	1	2	3	4	...	$n$
Cash flow	1	3	$3^2$	$3^3$	$3^4$	...	$3^n$

- (b) Derive the present value  $S$  of the following (infinite) cash flow in terms of  $r$  and/or  $x$ :

Period	0	1	2	3	4	...	$n$
Cash flow	1	2	3	4	5	...	$n+1$

7. (Callable bond) The Z Corporation issues a 10%, 20-year bond at a time when yields are 10%. The bond has a call provision that allows the corporation to force a bond holder to redeem his or her bond at face value plus 5%. After 5 years the corporation finds that exercise of this call provision is advantageous. What can you deduce about the yield at that time? (Assume one coupon payment per year.)
8. (The biweekly mortgage  $\oplus$ ) Here is a proposal that has been advanced as a way for homeowners to save thousands of dollars on mortgage payments: pay biweekly instead of monthly. Specifically, if monthly payments are  $x$ , it is suggested that one instead pay  $x/2$  every two weeks (for a total of 26 payments per year). This will pay down the mortgage faster, saving interest. The savings are surprisingly dramatic for this seemingly minor modification—often cutting the total interest payment by over one-third. Assume a loan amount of \$100,000 for 30 years at 10% interest, compounded monthly.
- (a) Under a monthly payment program, what are the monthly payments and the total interest paid over the course of the 30 years?
  - (b) Using the biweekly program, when will the loan be completely repaid, and what are the savings in total interest paid over the monthly program? (You may assume biweekly compounding for this part.)
9. (Annual worth) One advantage of the annual worth method is that it simplifies the comparison of investment projects that are repetitive but have different cycle times. Consider the automobile purchase problem of Example 2.7. Find the annual worths of the two (single-cycle) options, and determine directly which is preferable.
10. (Variable-rate mortgage  $\oplus$ ) The Smith family just took out a variable-rate mortgage on their new home. The mortgage value is \$100,000, the term is 30 years, and initially the interest rate is 8%. The interest rate is guaranteed for 5 years, after which time the rate will be adjusted according to prevailing rates. The new rate can be applied to their loan either by changing the payment amount or by changing the length of the mortgage.
- (a) What is the original yearly mortgage payment? (Assume payments are yearly.)
  - (b) What will be the mortgage balance after 5 years?
  - (c) If the interest rate on the mortgage changes to 9% after 5 years, what will be the new yearly payment that keeps the termination time the same?
  - (d) Under the interest change in (c), what will be the new term if the payments remain the same?
11. (Bond market) There are three bonds in the market as follows:
1. A bond with 4% coupon rate (paid annually), 10 years to maturity, and \$1,000 face value
  2. A bond with 4% plus current (short) rate (paid annually), 10 years to maturity, and \$1,000 face value
  3. A bond with 8% minus current (short) rate (paid annually), 10 years to maturity, and \$1,000 face value
- The prices of the bonds are \$950, \$1,100, and \$900, respectively.
- (a) Derive the price of a zero-coupon bond with 10 years to maturity and \$1,000 face value.
  - (b) Derive the price of a floating-rate bond (coupon paid annually) with 10 years to maturity and \$1,000 face value.
12. (Inflation-adjusted bonds) In 1997, the U.S. Treasury issued “Treasury inflation-protected securities” (**TIPS**). These are fixed-income securities that are inflation indexed to

protect their value against inflation. Like conventional bonds they have a fixed coupon rate and maturity date, but the face value is periodically adjusted for inflation by multiplying the original face value by the ratio of the Consumer Price Index (CPI) at the current date to the CPI at the original issue date. Each semiannual coupon payment is the inflation-adjusted face value times the fixed coupon rate. At maturity, the bondholder receives the maximum of the inflation-adjusted face value or the original face value. Hence, if deflation occurs, the bondholder is guaranteed not to lose on the face value.

You observe the following prices of two 10-year inflation-adjusted bonds:

$$\text{Bond 1: } P_1 = 77.92, C_1 = 3\%, F_1 = 100$$

$$\text{Bond 2: } P_2 = 100.00, C_2 = 6\%, F_2 = 100,$$

where  $P$  is the price,  $C$  is the coupon rate, and  $F$  is the original face value.

- (a) Compute the price of a theoretical 10-year inflation-adjusted zero-coupon bond with original face value of 100. (*Note:* Ignore taxes.)
- (b) Does your answer in (a) depend on the rate of inflation? Justify your answer.

**13. (Bond price)** An 8% bond with 18 years to maturity has a yield of 9%. What is the price of this bond?

**14. (Duration)** Find the price and duration of a 10-year, 8% bond that is trading at a yield of 10%.

**15. (Annuity duration  $\diamond$ )** Find the duration  $D$  and the modified duration  $D_M$  of a perpetual annuity that pays an amount  $A$  at the beginning of each year, with the first such payment being 1 year from now. Assume a constant interest rate  $r$  compounded yearly. [*Hint:* It is not necessary to evaluate any new summations.]

**16. (Bond selection)** Consider the four bonds having annual payments as shown in Table 3.9. They are traded to produce a 15% yield.

- (a) Determine the price of each bond.
- (b) Determine the duration of each bond (*not* the modified duration).
- (c) Which bond is most sensitive to a change in yield?
- (d) Suppose you owe \$2,000 at the end of 2 years. Concern about interest rate risk suggests that a portfolio consisting of the bonds and the obligation should be immunized. If  $V_A$ ,  $V_B$ ,  $V_C$ , and  $V_D$  are the total values of bonds purchased of types A, B, C, and D, respectively, what are the necessary constraints to implement the immunization? [*Hint:* There are two equations. (Do not solve.)]

**TABLE 3.9**  
**END-OF-YEAR PAYMENTS**

	Bond A	Bond B	Bond C	Bond D
Year 1	100	50	0	$0 + 1,000$
Year 2	100	50	0	0
Year 3	$100 + 1,000$	$50 + 1,000$	$0 + 1,000$	0

- (e) In order to immunize the portfolio, you decide to use bond C and one other bond. Which other bond should you choose? Find the amounts (in total value) of each of these to purchase.
- (f) You decided in (e) to use bond C in the immunization. Would other choices, including perhaps a combination of bonds, lead to lower total cost?

- 17.** (Continuous compounding  $\diamond$ ) Under continuous compounding the Macaulay duration becomes

$$D = \frac{\sum_{k=0}^n t_k e^{-\lambda t_k} c_k}{P},$$

where  $\lambda$  is the yield and

$$P = \sum_{k=0}^n e^{-\lambda t_k} c_k.$$

Find  $dP/d\lambda$  in terms of  $D$  and  $P$ .

- 18.** (Duration limit) Show that the limiting value of duration as maturity is increased to infinity is

$$D \rightarrow \frac{1 + (\lambda/m)}{\lambda}.$$

For the bonds in Table 3.6 (where  $\lambda = .05$  and  $m = 2$ ) we obtain  $D \rightarrow 20.5$ . Note that for large  $\lambda$  this limiting value approaches  $1/m$ , and hence the duration for large yields tends to be relatively short.

- 19.** (Convexity value) Find the convexity of a zero-coupon bond maturing at time  $T$  under continuous compounding (that is, when  $m \rightarrow \infty$ ).

- 20.** (Convexity theorem  $\diamond$ ) Suppose that an obligation occurring at a single time period is immunized against interest rate changes with bonds that have only nonnegative cash flows (as in the X Corporation example). Let  $P(\lambda)$  be the value of the resulting portfolio, including the obligation, when the interest rate is  $r + \lambda$  and  $r$  is the current interest rate. By construction  $P(0) = 0$  and  $P'(0) = 0$ . In this exercise we show that  $P(0)$  is a local minimum; that is,  $P''(0) \geq 0$ . (This property is exhibited by Example 3.10.)

Assume a yearly compounding convention. The discount factor for time  $t$  is  $d_t(\lambda) = (1 + r + \lambda)^{-t}$ . Let  $d_t = d_t(0)$ . For convenience assume that the obligation has magnitude 1 and is due at time  $\bar{t}$ . The conditions for immunization are then

$$\begin{aligned} P(0) &= \sum_t c_t d_t - d_{\bar{t}} = 0 \\ P'(0)(1+r) &= \sum_t t c_t d_t - \bar{t} d_{\bar{t}} = 0. \end{aligned}$$

- (a) Show that for all values of  $\alpha$  and  $\beta$  there holds

$$P''(0)(1+r)^2 = \sum_t (t^2 + \alpha t + \beta) c_t d_t - (\bar{t}^2 + \alpha \bar{t} + \beta) d_{\bar{t}}.$$

- (b) Show that  $\alpha$  and  $\beta$  can be selected so that the function  $t^2 + \alpha t + \beta$  has a minimum at  $\bar{t}$  and has a value of 1 there. Use these values to conclude that  $P''(0) \geq 0$ .

## References

The money market is vast and consists of numerous financial instruments and institutions. Detailed descriptions are available from many sources. Some good starting points are [1–5]. For comprehensive treatments of yield curve analysis, see [5–7]. The concept of *duration* was invented by Macaulay and by Redington; see [8, 9]. For history and details on the elaboration of this concept into a full methodology for immunization, see [10–13]. The result of Exercise 20 is a version of the Fisher-Weil theorem [13].

1. Cook, T. Q., and T. D. Rowe (1986), *Instruments of the Money Market*, Federal Reserve Bank, Richmond, VA.
2. Fabozzi, F. J., and F. Modigliani (2008), *Capital Markets: Institutions and Instruments*, 4th ed., Prentice Hall, Englewood Cliffs, NJ.
3. *Handbook of U.S. Government and Federal Agency Securities and Related Money Market Instruments*, "The Pink Book," 34th ed. (1990), The First Boston Corporation, Boston.
4. Homer, S., and M. Liebowitz (1972), *Inside the Yield Book: New Tools for Bond Market Strategy*, 4th ed., Prentice Hall, Englewood Cliffs, NJ.
5. Veronesi, P. (2010), *Fixed-Income Securities*, John Wiley & Sons, Hoboken, NJ.
6. Fabozzi, F. J. (2001), *Fixed-Income Securities*, 2nd ed., John Wiley & Sons, New York.
7. Van Home, J. C. (2000), *Financial Market Rates and Flows*, 6th ed., Prentice Hall, Englewood Cliffs, NJ.
8. Macaulay, F. R. (1938), *Some Theoretical Problems Suggested by the Movement of Interest Rates, Bond Yield, and Stock Prices in the United States since 1856*, National Bureau of Economic Research, New York.
9. Redington, F. M. (October 1971), "Review of the Principles of Life-Office Valuations," *Journal of the Institute of Actuaries*, **78**, no. 3, 286–315.
10. Bierwag, G. O., and G. G. Kaufman (July 1977), "Coping with the Risk of Interest-Rate Fluctuations: A Note," *Journal of Business*, **50**, no. 3, 364–370.
11. Bierwag, G. O., G. G. Kaufman, and A. Toebs (July–August 1983), "Duration: Its Development and Use in Bond Portfolio Management," *Financial Analysts Journal*, **39**, no. 4, 15–35.
12. Bierwag, G. O. (1987), *Duration Analysis*, Ballinger Publishing, Cambridge, MA.
13. Fisher, L., and R. L. Weil (1971), "Coping with the Risk of Interest-Rate Fluctuations: Returns to Bondholders from Naive and Optimal Strategies," *Journal of Business*, **44**, 408–431.

# 4

## THE TERM STRUCTURE OF INTEREST RATES

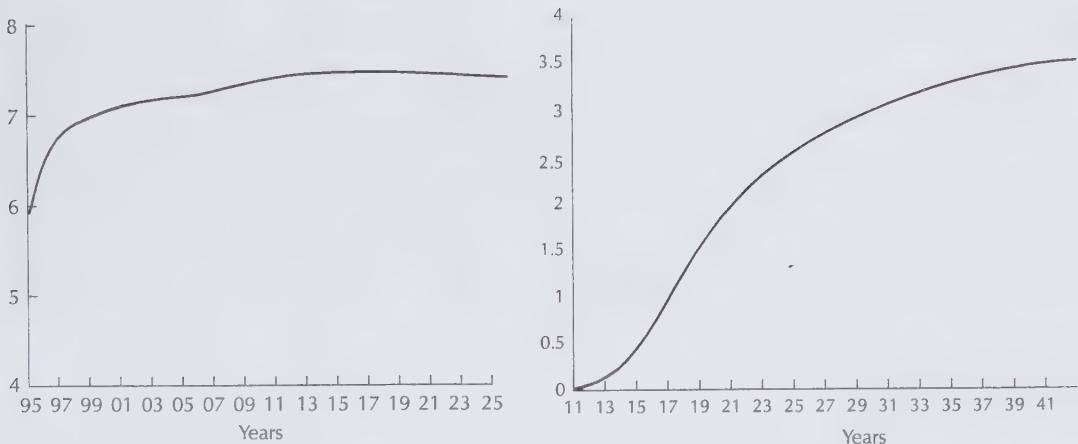
A richer theory of interest rates is explored in this chapter, as compared to that in previous chapters. The enriched theory allows for a whole family of interest rates at any one time—a different rate for each maturity time—providing a clearer understanding of the interest rate market and a foundation for more sophisticated investment analysis techniques.

### 4.1 The Yield Curve

The yield to maturity of any bond is strongly tied to general conditions in the fixed-income securities market. All yields tend to move together in this market. However, all bond yields are not exactly the same.

The variation in yields across bonds is explained in part by the fact that bonds have various quality ratings. A strong AAA-rated bond is likely to cost more (and hence have lower yield) than a bond with an identical promised income stream but having a B-quality rating. It is only natural that high quality is more expensive than low quality. However, quality alone does not fully explain the observed variations in bond yields.

Another factor that partially explains the differences in the yields of various bonds is the time to maturity. The situation can be represented by a **yield curve**, which is a plot of yield versus maturity for bonds of similar quality. As a general rule, “long” bonds with very distant maturity dates tend to offer greater yields than “short” bonds of the same quality. Such curves are shown in Figure 4.1. The curve on the left increases



**FIGURE 4.1 Yield curve.** Yields are plotted as a function of maturity date. The curve on the left, from 1995, is typical. The curve on the right, from 2011, is unusual. Source: *Treasury Bulletin*, June 1995, September 2011.

monotonically with maturity. Such a yield curve is said to be “normally shaped,” reflecting the fact that long-term bonds are usually less desirable than short-term bonds. If there are regions where the yield curve decreases as maturity increases, the curve is said to be **inverted**. The inverted shape tends to occur when short-term rates increase rapidly and investors believe that the rise is temporary, so long-term rates remain near their previous levels. There are times when the curve is fully inverted, decreasing everywhere. If the curve decreases only in a small region of maturities, the curve is said to be *partially* inverted. The yield curve on the right of Figure 4.1 is increasing monotonically everywhere, so it is not inverted, but it has a region at short maturities where rates are near zero and the curve is convex, as if the curve was “squashed.” This can be the result of the Federal Reserve’s actively suppressing short-term interest rates. The yield curve often undulates over time, somewhat like a branch in the wind.

When studying a particular bond, it is useful to determine its yield and maturity date and place it as a point on the yield curve for bonds in its risk class. This will give a general indication of how it is priced relative to the overall market. If it is far from the curve, there is probably a reason, related to special situations or special features (such as call features of the bond or news affecting the potential solvency of the issuer).

The yield curve is helpful, but because it is a bit arbitrary, it does not provide a completely satisfactory explanation of yield differences. Why, for example, should the maturity date be used as the horizontal axis of the curve rather than, say, duration? A more basic theory is required, and such a theory is introduced in the next section.

## 4.2 The Term Structure

Term structure theory puts aside the notion of yield and instead focuses on pure interest rates. The theory is based on the observation that, in general, the interest rate charged (or paid) for money depends on the length of time that the money is held. Your local bank, for example, is likely to offer you a higher rate of interest for deposits committed for 3 years than for demand deposits (which can be withdrawn at any time). This basic fact, that the interest rate charged depends on the length of time that the funds are held, is the basis of term structure theory. This chapter works out the details and implications of that fact.

### Spot Rates

**Spot rates** are the basic interest rates defining the term structure. The spot rate  $s_t$  is the rate of interest, expressed in yearly terms, charged for money held from the present time ( $t = 0$ ) until time  $t$ . Both the interest and the original principal are paid at time  $t$ . Hence, in particular,  $s_1$  is the 1-year interest rate; that is, it is the rate paid for money held 1 year. Similarly, the rate  $s_2$  represents the rate that is paid for money held 2 years; however, it is expressed on an annualized basis. Thus if your bank promises to pay a rate of  $s_2$  for a 2-year deposit of an amount  $A$  compounded yearly, it will actually repay  $(1 + s_2)^2 A$  at the end of 2 years; your money grows by a factor of  $(1 + s_2)^2$ .

The definition of spot rates implicitly assumes a compounding convention, and this convention might vary with the purpose at hand. The preceding discussion assumed a 1-year compounding convention. It is common to use  $m$  periods per year, or continuous compounding, as well. In all cases the rates are usually still quoted as yearly rates. For completeness, we list the various possibilities:

- (a) **Yearly** Under the yearly compounding convention, the spot rate  $s_t$  is defined such that

$$(1 + s_t)^t$$

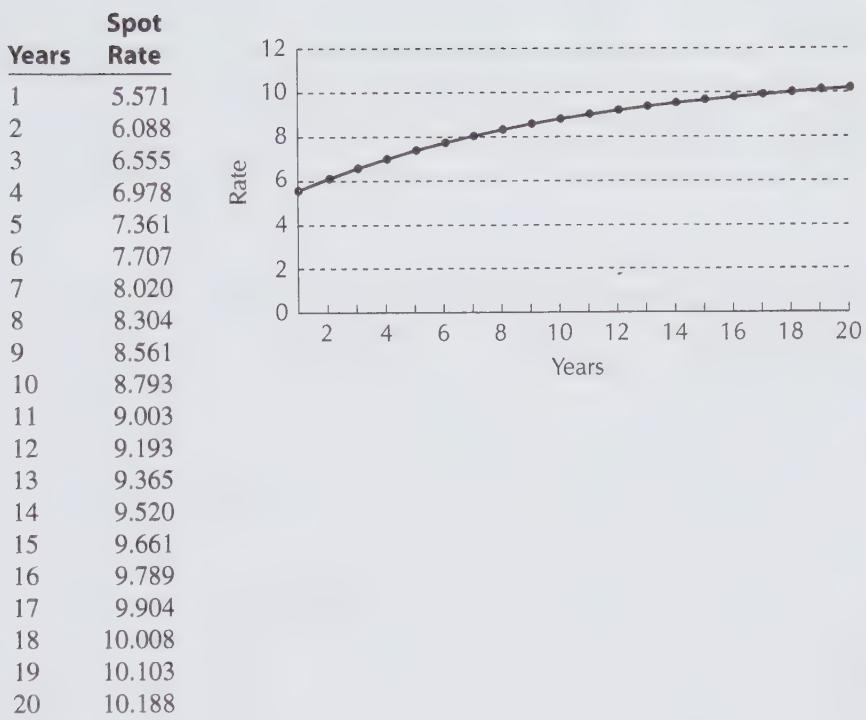
is the factor by which a deposit held  $t$  years will grow. (Here,  $t$  must be an integer, or an adjustment must be made.)

- (b)  **$m$  periods per year** Under a convention of compounding  $m$  periods per year, the spot rate  $s_t$  is defined so that

$$(1 + s_t/m)^{mt}$$

is the corresponding factor. (Here,  $mt$  must be an integer, so  $t$  must be an integral multiple of  $1/m$ .)

- (c) **Continuous** Under a continuous compounding convention, the spot rate  $s_t$  is defined so that  $e^{s_t t}$  is the corresponding growth factor. This formula applies directly to all values of  $t$ .



**FIGURE 4.2** Spot rate curve. The yearly rate of interest depends on the length of time funds are held.

For theoretical purposes, continuous compounding is “neater” since the formulas apply without change to all values of  $t$ . The other methods require an adjustment for values of  $t$  between compounding dates. However, the yearly compounding convention is the most convenient, and it is the convention mainly used in this chapter.

Spot rates can, in theory, be measured by recording the yields of zero-coupon bonds. (In order to eliminate the influence of default risk, it would be best to consider only Treasury securities for this purpose.) Since a zero-coupon bond promises to pay a fixed amount at a fixed date in the future, the ratio of the payment amount to the current price defines the spot rate for the maturity date of the bond. By this measurement process we can develop a spot rate curve, which is analogous to the yield curve. Such a curve and a chart of the corresponding data are shown in Figure 4.2.

## Discount Factors and Present Value

Once the spot rates have been determined, it is natural to define the corresponding **discount factors**  $d_t$  for each time point. These are the factors by which future

cash flows must be multiplied to obtain an equivalent present value. For the various compounding conventions, they are defined as follows:

(a) **Yearly** For yearly compounding,

$$d_k = \frac{1}{(1 + s_k)^k}.$$

(b)  **$m$  periods per year** For compounding  $m$  periods per year,

$$d_k = \frac{1}{(1 + s_k/m)^{mk}}.$$

(a) **Continuous** For continuous compounding,

$$d_t = e^{-s_t t}.$$

The discount factors transform future cash flows directly into an equivalent present value. Hence given any cash flow stream  $(x_0, x_1, x_2, \dots, x_n)$ , the present value, relative to the prevailing spot rates, is

$$PV = x_0 + d_1 x_1 + d_2 x_2 + \dots + d_n x_n.$$

The discount factor  $d_k$  acts like a *price* for cash received at time  $k$ . We determine the value of a stream by adding up “price times quantity” for all the cash components of the stream.

**Example 4.1 (Price of a 10-year bond)** Using the spot rate curve of Figure 4.2, let us find the value of an 8% bond maturing in 10 years.

Normally, for bonds we would use the rates and formulas for 6-month compounding; but for this example let us assume that coupons are paid only at the end of each year, starting a year from now, and that 1-year compounding is consistent with our general evaluation method. We write the cash flows together with the discount factors, take their products, and then sum, as shown in Table 4.1. The value of the bond is found to be 97.34.

**TABLE 4.1  
BOND EVALUATION**

Year	1	2	3	4	5	6	7	8	9	10	Total PV
Discount	.947	.889	.827	.764	.701	.641	.583	.528	.477	.431	
Cash flow	8	8	8	8	8	8	8	8	8	108	
PV	7.58	7.11	6.61	6.11	5.61	5.12	4.66	4.22	3.82	46.50	97.34

*Each cash flow is discounted by the discount factor for its time.*

**Example 4.2 (Simplico gold mine)** Consider the lease of the Simplico gold mine discussed in Chapter 2, Example 2.6, but now let us assume that interest rates follow the term structure pattern of Figure 4.2. We shall find the present value of the lease.

The cash flow stream is identical to that of the earlier example; namely, \$2M each year for 10 years. The present value is therefore just the sum of the first 10 discount figures multiplied by \$2M, for a total of \$13.58M.

## Determining the Spot Rate

The obvious way to determine a spot rate curve is to find the prices of a series of zero-coupon bonds with various maturity dates. Unfortunately the set of available zero-coupon bonds is typically rather sparse, and, indeed, until recently there were essentially no “zeros” available with long maturities. Thus it is not always practical to determine a complete set of spot rates this way. However, the existence of zero-coupon bonds is not necessary for the concept of spot rates to be useful, nor are they needed as data to determine the spot rate value.

The spot rate curve can be determined from the prices of coupon-bearing bonds by beginning with short maturities and working forward toward longer maturities. We illustrate the process for the 1-year compounding convention (and assuming coupons are paid only once a year). First determine  $s_1$  by direct observation of the 1-year interest rate—as determined, for example, by the 1-year Treasury bill rate. Next consider a 2-year bond. Suppose that bond has price  $P$ , makes coupon payments of amount  $C$  at the end of both years, and has a face value  $F$ . The price should equal the discounted value of the cash flow stream, so we can write

$$P = \frac{C}{1+s_1} + \frac{C+F}{(1+s_2)^2}.$$

Since  $s_1$  is already known, we can solve this equation for  $s_2$ . Working forward this way, by next considering 3-year bonds, then 4-year bonds, and so forth, we can determine  $s_3, s_4, \dots$ , step by step.

Spot rates can also be determined by a subtraction process. Two bonds of different coupon rates but identical maturity dates can be used to construct the equivalent of a zero-coupon bond. The following example illustrates the method.

**Example 4.3 (Construction of a zero)** Bond A is a 10-year bond with a 10% coupon. Its price is  $P_A = 98.72$ . Bond B is a 10-year bond with an 8% coupon. Its price is  $P_B = 85.89$ . Both bonds have the same face value, normalized to 100.

Consider a portfolio with  $-0.8$  unit of bond A and 1 unit of bond B. This portfolio will have a face value of 20 and a price of  $P = P_B - 0.8P_A = 6.914$ . The coupon payments cancel, so this is a zero-coupon portfolio. The 10-year spot rate  $s_{10}$  must satisfy  $(1+s_{10})^{10}P = 20$ . Thus  $s_{10} = 11.2\%$ .

In practice, since spot rates are an idealization and the spot rates implied by different bonds may differ slightly from one another, it is advisable to modify these

procedures to incorporate an averaging method when estimating the spot rates. (See Exercise 4.)

## 4.3 Forward Rates

An elegant and useful concept emerges directly from the definition of spot rates; namely, the concept of forward rates. **Forward rates** are interest rates for money to be borrowed between two dates in the future, *but under terms agreed upon today*.

It is easiest to explain the concept for a 2-year situation. Suppose that  $s_1$  and  $s_2$  are known. If we leave \$1 in a 2-year account it will, by definition, grow to  $\$(1 + s_2)^2$  at the end of the 2 years. Alternatively, we might place the \$1 in a 1-year account and simultaneously make arrangements that the proceeds, which will be  $\$(1 + s_1)$ , will be lent for 1 year starting a year from now. That loan will accrue interest at a prearranged rate (agreed upon now) of say  $f$ . The rate  $f$  is the forward rate for money to be lent in this way. The final amount of money we receive at the end of 2 years under this compound plan is  $\$(1 + s_1)(1 + f)$ .

We now invoke the comparison principle. We have two alternative methods for investing \$1 for 2 years. The first returns  $(1 + s_2)^2$  and the second returns  $(1 + s_1)(1 + f)$ .

These two should be equal, since both are available.<sup>1</sup> Thus we have

$$(1 + s_2)^2 = (1 + s_1)(1 + f)$$

or

$$f = \frac{(1 + s_2)^2}{1 + s_1} - 1.$$

Hence the forward rate is determined by the two spot rates.

We can justify the use of the comparison principle here through an **arbitrage argument**. If these two methods of investing money did not return the same amount, then there would be an opportunity to make arbitrage profits—defined to be a chance of riskless profit. In the preceding example, if  $(1 + s_1)(1 + f) > (1 + s_2)^2$ , meaning that the second method of investment returned more than the first, then an arbitrageur could reverse the first plan (by *borrowing* for 2 years) and then carry out the second plan by investing the money that was borrowed. This arbitrageur would have zero net investment because he or she used only borrowed capital, but after repaying the loan the arbitrageur would have a profit factor of  $(1 + s_1)(1 + f) - (1 + s_2)^2 > 0$ . This arbitrage scheme could be carried out at any magnitude, and hence, in theory, the arbitrageur could make very large sums of money from no initial capital. We assume that it is not possible to implement this scheme in the market because potential arbitrageurs are always on the lookout for such discrepancies. If a slight discrepancy does arise, they take advantage of it, and this action tends to close the gap in rates.

---

<sup>1</sup> Forward contracts of this type are actually implemented by the use of futures contracts on Treasury securities, as explained in Chapter 12. They are highly liquid, so forwards of this type are obtained easily.

If the inequality were in the other direction, the arbitrageur could just reverse the procedure. Thus equality must hold.

The arbitrage argument assumes that there are no transaction costs—either real costs such as brokerage fees or opportunity costs related to the time and effort of finding the discrepancy and arranging for the trades. The argument also assumes that the borrowing and lending rates are identical. If there were transaction costs or unequal rates, there could be a slight “wedge” between the 2-year rates associated with the two alternative strategies. However, in practice the transaction cost associated with a highly liquid security such as a U.S. Treasury is a very small fraction of the security’s total cost, especially if large amounts are involved; and borrowing and lending rates are quite close, again if large amounts of capital are involved. So although the arbitrage argument represents an idealization, it is in practice a reasonable approximation.

The comparison principle can be used to argue that the two overall rates must be equal even in the absence of arbitrageurs. If there were a difference in rates, then investors seeking to loan money for 2 years would choose the best alternative—and so would borrowers. Market forces would tend to equalize the rates.

**Example 4.4** Suppose that the spot rates for 1 and 2 years are, respectively,  $s_1 = 7\%$  and  $s_2 = 8\%$ . We then find that the forward rate is  $f = (1.08)^2 / 1.07 - 1 = .0901 = 9.01\%$ . Hence the 2-year 8% rate can be obtained either as a direct 2-year investment, or by investing for 1 year at 7% followed by a second year at 9.01%.

This discussion can be generalized to define other forward rates between different time periods. The rate  $f$  used earlier is more completely labeled as  $f_{1,2}$  because it is the forward rate between years 1 and 2. In general we use the following:

**Forward rate definition** *The forward rate between times  $t_1$  and  $t_2$  with  $t_1 < t_2$  is denoted by  $f_{t_1,t_2}$ . It is the rate of interest charged for borrowing money at time  $t_1$  which is to be repaid (with interest) at time  $t_2$ .*

In general, forward rates are expressed on an annualized basis, like other interest rates, unless another basis is explicitly specified.

As mentioned before, in the market there could be more than one rate for any particular forward period. For example, the forward rate for borrowing may differ from that for lending. Thus when discussing market rates one must be specific. However, in theoretical discussions the definition of forward rates is based on an underlying set of spot rates (which themselves generally represent idealizations or averages of market conditions). These calculated forward rates are often termed **implied forward rates** to distinguish them from **market forward rates**.

The implied forward rates are found by extending the logic given earlier for assigning the value  $f_{1,2}$ . If we use 1-year compounding, the basic forward rates are defined between various yearly periods. They are defined to satisfy the following equation (for  $i < j$ ):

$$(1 + s_j)^j = (1 + s_i)^i (1 + f_{i,j})^{j-i}.$$

The left side of this equation is the factor by which money grows if it is directly invested for  $j$  years. This amount is determined by the spot rate  $s_j$ . The right side of the equation is the factor by which money grows if it is invested first for  $i$  years and then in a forward contract (arranged now) between years  $i$  and  $j$ . The term  $(1 + f_{i,j})$  is raised to the  $(j - i)$ th power because the forward rate is expressed in yearly terms.

The extension to other compounding conventions is straightforward. For completeness, the formulas for forward rates (expressed as yearly rates) under various compounding conventions are listed here:

**Forward rate formulas** *The implied forward rate between times  $t_1$  and  $t_2 > t_1$  is the rate of interest between those times that is consistent with a given spot rate curve. Under various compounding conventions the forward rates are specified as follows:*

(a) **Yearly** *For yearly compounding, the forward rates satisfy, for  $j > i$ ,*

$$(1 + s_j)^j = (1 + s_i)^i (1 + f_{i,j})^{j-i}$$

*Hence,*

$$f_{i,j} = \left[ \frac{(1 + s_j)^j}{(1 + s_i)^i} \right]^{1/(j-i)} - 1.$$

(b)  **$m$  periods per year** *For  $m$  period-per-year compounding, the forward rates satisfy, for  $j > i$ , expressed in periods,*

$$(1 + s_i/m)^j = (1 + s_i/m)^i (1 + f_{i,j}/m)^{(j-i)}.$$

*Hence,*

$$f_{i,j} = m \left[ \frac{(1 + s_i/m)^j}{(1 + s_i/m)^i} \right]^{1/(j-i)} - m.$$

(c) **Continuous** *For continuous compounding, the forward rates  $f_{t_1,t_2}$  are derived for all  $t_1$  and  $t_2$ , with  $t_2 > t_1$ , and satisfy*

$$e^{s_{t_2} t_2} = e^{s_{t_1} t_1} e^{f_{t_1,t_2} (t_2 - t_1)}.$$

*Hence,*

$$f_{t_1,t_2} = \frac{s_{t_2} t_2 - s_{t_1} t_1}{t_2 - t_1}.$$

Note again that continuous compounding produces the simplest formula. As a further convention, it is useful to define spot rates, discount factors, and forward rates when one of the time points is zero, representing current time. Hence we define  $s_{t_0} = 0$  and correspondingly  $d_{t_0} = 1$ , where  $t_0$  is the current time. (Alternatively we

write  $s_0 = 0$  and  $d_0 = 1$  when denoting time by period integers.) For forward rates, we write similarly  $f_{t_0,t_1} = s_{t_1}$ . The forward rates from time zero are the corresponding spot rates.

There are a large number of forward rates associated with a spot rate curve. In fact, if there are  $n$  periods, there are  $n$  spot rates (excluding  $s_0$ ); and there are  $n(n + 1)/2$  forward rates (including the basic spot rates.) However, all these forward rates are derived from the  $n$  underlying spot rates.

The forward rates are introduced partly because they represent rates of actual transactions. Forward contracts do in fact serve a very important hedging role, and their use in this manner is discussed further in Chapter 12. They are introduced here, however, mainly because they are important for the full development of the term structure theory. They are used briefly in the next section and then extensively in the section following that.

## 4.4 Term Structure Explanations

The yield curve can be observed, at least roughly, by looking at a series of bond quotes in the financial press. The curve is almost never flat but, rather, it usually slopes gradually upward as maturity increases. The spot rate curve has similar characteristics. Typically it, too, slopes rapidly upward at short maturities and continues to slope upward, but more gradually as maturities lengthen. It is natural to ask if there is a simple explanation for this typical shape. Why is the curve not just flat at a common interest rate?

There are three standard explanations (or “theories”) for the term structure, each of which provides some important insight. We outline them briefly in this section.

### Expectations Theory

The first explanation is that spot rates are determined by expectations of what rates will be in the future. To visualize this process, suppose that, as is usually the case, the spot rate curve slopes upward, with rates increasing for longer maturities. The 2-year rate is greater than the 1-year rate. It is argued that this is so because the market (that is, the collective of all people who trade in the interest rate market) believes that the 1-year rate will most likely go up next year. (This belief may, for example, be because most people believe inflation will rise, and thus to maintain the same real rate of interest, the nominal rate must increase.) This majority belief that the interest rate will rise translates into a market *expectation*. An expectation is only an average guess; it is not definite information—for no one knows for sure what will happen next year—but people on average assume, according to this explanation, that the rate will increase.

This argument is made more concrete by expressing the expectations in terms of forward rates. This more precise formulation is the **expectations hypothesis**. To outline this hypothesis, consider the forward rate  $f_{1,2}$ , which is the implied rate for

money loaned for 1 year, a year from now. According to the expectations hypothesis, this forward rate is *exactly* equal to the market expectation of what the 1-year spot rate will be next year. Thus the expectation can be inferred from existing rates.

Earlier we considered a situation where  $s_1 = 7\%$  and  $s_2 = 8\%$ . We found that the implied forward rate was  $f_{1,2} = 9.01\%$ . According to the unbiased expectations hypothesis, this value of 9.01% is the market's expected value of next year's 1-year spot rate  $s'_1$ .

The same argument applies to the other rates as well. As additional spot rates are considered, they define corresponding forward rates for next year. Specifically,  $s_1, s_2$ , and  $s_3$  together determine the forward rates  $f_{1,2}$  and  $f_{1,3}$ . The second of these is the forward rate for borrowing money for 2 years, starting next year. This rate is assumed to be equal to the current expectation of what the 2-year spot rate  $s'_2$  will be next year. In general, then, the current spot rate curve leads to a set of forward rates  $f_{1,2}, f_{1,3}, \dots, f_{1,n}$ , which define the expected spot rate curve  $s'_1, s'_2, \dots, s'_{n-1}$  for next year. The expectations are inherent in the current spot rate structure.

There are two ways of looking at this construction. One way is that the current spot rate curve implies an expectation about what the spot rate curve will be next year. The other way is to turn this first view around and say that the expectation of next year's curve determines what the current spot rate curve must be. Both views are intertwined; expectations about future rates are part of today's market and influence today's rates.

This theory or hypothesis is a nice explanation of the spot rate curve, even though it has some important weaknesses. The primary weakness is that, according to this explanation, the market expects rates to increase whenever the spot rate curve slopes upward; and this is practically all the time. Thus the expectations cannot be right even on average, since rates do not go up as often as expectations would imply. Nevertheless, the expectations explanation is plausible, although the expectations may themselves be skewed.

The expectations explanation of the term structure can be regarded as being (loosely) based on the comparison principle. To see this, consider again the 2-year situation. An investor can invest either in a 2-year instrument or in a 1-year instrument followed by another 1-year investment. The follow-on investment can also be carried out two ways. It can be arranged currently through a forward contract at rate  $f_{1,2}$ , or it can simply be "rolled over" by reinvesting the following year at the then prevailing 1-year rate. A wise investor would compare the two alternatives. If the investor expects that next year's 1-year rate will equal the current value of  $f_{1,2}$ , then he or she will be indifferent between these two alternatives. Indeed, the fact that both are viable implies that they must seem (approximately) equal.

## Liquidity Preference

The liquidity preference explanation asserts that investors usually prefer short-term fixed income securities over long-term securities. The simplest justification for this assertion is that investors do not like to tie up capital in long-term securities, since those funds may be needed before the maturity date. Investors prefer their funds to be

**liquid** rather than tied up. However, the term *liquidity* is used in a slightly nonstandard way in this argument. There are large active markets for bonds of major corporations and of the Treasury, so it is easy to sell any such bonds one might hold. Short-term and long-term bonds of this type are equally liquid.

Liquidity is used in this explanation of the term structure shape instead to express the fact that most investors prefer short-term bonds to long-term bonds. The reason for this preference is that investors anticipate that they may need to sell their bonds soon, and they recognize that long-term bonds are more sensitive to interest rate changes than are short-term bonds. Hence an investor who may need funds in a year or so will be reluctant to place these funds in long-term bonds because of the relatively high near-term risk associated with such bonds. To lessen risk, such an investor prefers short-term investments. Hence to induce investors into long-term instruments, better rates must be offered for long bonds. For this reason, according to the theory, the spot rate curve rises.

## Market Segmentation

The market segmentation explanation of the term structure argues that the market for fixed-income securities is segmented by maturity dates. This argument assumes that investors have a good idea of the maturity date that they desire, based on their projected need for future funds or their risk preference. The argument concludes that the group of investors competing for long-term bonds is different from the group competing for short-term bonds. Hence there need be no relation between the prices (defined by interest rates) of these two types of instruments; short and long rates can move around rather independently. Taken to an extreme, this viewpoint suggests that all points on the spot rate curve are mutually independent. Each is determined by the forces of supply and demand in its own market.

A moderated version of this explanation is that, although the market is basically segmented, individual investors are willing to shift segments if the rates in an adjacent segment are substantially more attractive than those of the main target segment. Adjacent rates cannot become grossly out of line with each other. Hence the spot rate curve must indeed be a curve rather than a jumble of disjointed numbers, but this curve can bend in various ways, depending on market forces.

## Discussion

Certainly each of the foregoing explanations embodies an element of truth. The whole truth is probably some combination of them all.

The expectations theory is the most analytical of the three, in the sense that it offers concrete numerical values for expectations, and hence it can be tested. These tests show that it works reasonably well with a deviation that seems to be explained by liquidity preference. Hence expectations tempered by the risk considerations of liquidity preference seem to offer a good straightforward explanation.

## 4.5 Expectations Dynamics

The concept of market expectations introduced in the previous section as an explanation for the shape of the spot rate curve can be developed into a useful tool in its own right. This tool can be used to form a plausible **forecast** of future interest rates.

### Spot Rate Forecasts

The basis of this method is to assume that the expectations implied by the current spot rate curve will actually be fulfilled. Under this assumption we can then predict next year's spot rate curve from the current one. This new curve implies yet another set of expectations for the following year. If we assume that these, too, are fulfilled, we can predict ahead once again. Going forward in this way, an entire future of spot rate curves can be predicted. Of course, it is understood that these predicted spot rate curves are based on the assumption that expectations will be fulfilled (and we recognize that this may not happen), but once made, the assumption does provide a logical forecast.

Let us work out some of the details. We begin with the current spot rate curve  $s_1, s_2, \dots, s_n$ , and we wish to estimate next year's spot rate curve  $s'_1, s'_2, \dots, s'_{n-1}$ . The current forward rate  $f_{1,j}$  can be regarded as the expectation of what the interest rate will be next year—measured from next year's current time to a time  $j - 1$  years ahead—in other words,  $f_{1,j}$  is next year's spot rate  $s'_{i-1}$ . Explicitly,<sup>2</sup>

$$s'_{j-1} = f_{1,j} = \left[ \frac{(1+s_j)^j}{1+s_1} \right]^{1/(j-1)} - 1 \quad (4.1)$$

for  $1 < j \leq n$ . This is the basic formula for updating a spot rate curve under the assumption that expectations are fulfilled. Starting with the current curve, we obtain an estimate of next year's curve.

We term this transformation **expectations dynamics**, since it gives an explicit characterization of the dynamics of the spot rate curve based on the expectations assumption. Other assumptions are certainly possible. For instance, we could assume that the spot rate curve will remain unchanged, or that it will shift upward by a fixed amount, and so forth; however, expectations dynamics has a nice logical appeal.

The expectations process can be carried out for another step to obtain the spot rate curve for the third year, and so forth. Note, however, that if the original curve has finite length, each succeeding curve is shorter by one term—and hence the curves eventually become quite short. This problem can be rectified by initially assuming a very long (or infinite) spot rate curve, or by adding a new  $s_n$  term each year. This latter approach would require an additional hypothesis.

---

<sup>2</sup> Recall that this formula for  $f_{1,j}$  was given in Section 4.3. It is derived from the relation  $(1 + s_j)^j = (1 + f_{1,j})^{j-1}(1 + s_1)$ .

**Example 4.5 (A simple forecast)** Let us take as given the spot rate curve shown in the first row of the table. The second row is then the forecast of next year's spot rate curve under expectations dynamics. This row is found using equation (4.1).

	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
Current	6.00	6.45	6.80	7.10	7.36	7.56	7.77
Forecast	6.90	7.20	7.47	7.70	7.88	8.06	

The first two entries in the second row were computed as follows:

$$f_{1,2} = \frac{(1.0645)^2}{1.06} - 1 = .069$$

$$f_{1,3} = \left[ \frac{(1.068)^3}{1.06} \right]^{1/2} - 1 = .072.$$

All future spot rate curves implied by an initial spot rate curve can be displayed by listing all of the forward rates associated with the initial spot rate curve. Such a list is shown in a triangular array:

$$\begin{array}{ccccccc} f_{0,1} & f_{0,2} & f_{0,3} & \dots & f_{0,n-2} & f_{0,n-1} & f_{0,n} \\ f_{1,2} & f_{1,3} & f_{1,4} & \dots & f_{1,n-1} & f_{1,n} & \\ f_{2,3} & f_{2,4} & f_{2,5} & \dots & f_{2,n} & & \\ \vdots & \vdots & & & & & \\ f_{n-2,n-1} & f_{n-2,n} & & & & & \\ f_{n-1,n} & & & & & & \end{array}$$

The first row of the array lists the forward rates from the initial time. These are identical to the spot rates themselves; that is,  $s_j = f_{0,j}$  for all  $j$  with  $0 < j \leq n$ . The next row lists the forward rates from time 1. These will be next year's spot rates according to expectations dynamics. The third row will be the spot rates for the third year, and so forth.

## Discount Factors

Another important concept is that of a **discount factor** between two times. The discount factors are, of course, fundamental quantities used in present value calculations.

It is useful to apply a double indexing system to the discount factors paralleling the system used for forward rates. Accordingly, the symbol  $d_{j,k}$  denotes the discount factor used to discount cash received at time  $k$  back to an equivalent amount of cash at time  $j$ . The normal, time zero, discount factors are  $d_1 = d_{0,1}, d_2 = d_{0,2}, \dots, d_n = d_{0,n}$ . The discount factors can be expressed in terms of the forward rates as

$$d_{j,k} = \left[ \frac{1}{1+f_{j,k}} \right]^{k-j}$$

The discount factors are related by a compounding rule: to discount from time  $k$  back to time  $i$ , one can first discount from time  $k$  back to an intermediate time  $j$  and then discount from  $j$  back to  $i$ . In other words,  $d_{i,k} = d_{i,j}d_{j,k}$  for  $i < j < k$ .

**Discount factor relation** *The discount factor between periods  $i$  and  $j$  is defined as*

$$d_{i,j} = \left[ \frac{1}{1+f_{i,j}} \right]^{j-i}$$

*These factors satisfy the compounding rule*

$$d_{i,k} = d_{i,j}d_{j,k}$$

*for  $i < j < k$ .*

## Short Rates

Short rates are the forward rates spanning a single time period. The short rate at time  $k$  is accordingly  $r_k = f_{k,k+1}$ ; that is, it is the forward rate from  $k$  to  $k+1$ . The short rates can be considered fundamental just as spot rates, for a complete set of short rates fully specifies a term structure.

The spot rate  $s_k$  is found from the short rates from the fact that interest earned from time zero to time  $k$  is identical to the interest that would be earned by rolling over an investment each year. Specifically,

$$(1+s_k)^k = (1+r_0)(1+r_1)\cdots(1+r_{k-1}).$$

The relation generalizes because all forward rates can be found from the short rates in a similar way. Specifically,

$$(1+f_{i,j})^{j-i} = (1+r_i)(1+r_{i+1})\cdots(1+r_{j-1}).$$

Hence the short rates form a convenient basis for generating all other rates.

The short rates are especially appealing in the context of expectations dynamics, because they do not change from year to year, whereas spot rates do. Given the initial short rates  $r_0, r_1, r_2, \dots, r_{n-1}$ , next year (under expectations dynamics) the short rates will be  $r_1, r_2, \dots, r_{n-1}$ . The short rate for a specific year does not change; however, that year is 1 year closer to the sliding current time. For example, if we are at the beginning of year 2030, the short rate  $r_4$  is the rate for the year beginning January 2034. A year later, in 2031, the new  $r_3$  will be the rate for the year 2034, and this short rate will be identical (under expectations dynamics) to the previous  $r_4$ .

An example of a complete set of forward rates, discount factors, and short rates is shown in Table 4.2. Here the rows represent the rates or factors for a given year: the top row of each array contains the initial rates or factors for 7 years forward. The forward rate array is, as discussed, identical to the spot rate array. Hence the basic spot rate curve is defined by the top line of the forward rate array. Everything else is derived from that single row. The discount factors for the current time are those

**TABLE 4.2**  
**FORWARD RATES, DISCOUNT FACTORS, AND SHORT RATES**

Forward rates							Short rates						
6.00	6.45	6.80	7.10	7.36	7.56	7.77	6.00	6.90	7.50	8.00	8.40	8.60	9.00
6.90	7.20	7.47	7.70	7.88	8.06		6.90	7.50	8.00	8.40	8.60	9.00	
7.50	7.75	7.97	8.12	8.30			7.50	8.00	8.40	8.60	9.00		
8.00	8.20	8.33	8.50				8.00	8.40	8.60	9.00			
8.40	8.50	8.67					8.40	8.60	9.00				
8.60	8.80						8.60	9.00					
9.00							9.00						
Discount factors													
.943	.883	.821	.760	.701	.646	.592							
.935	.870	.806	.743	.684	.628								
.930	.861	.795	.732	.671									
.926	.854	.787	.722										
.923	.849	.779											
.921	.845												
.917													

The original spot rate curve is defined by the top row of the forward rate array. All other terms are derived from this row.

listed in the top row of the discount factor array. These are the values used to find the present values of future cash flows. Note that successive rows of the short rate table are just shifted versions of the rows above. Short rates remain fixed in absolute time.

## Invariance Theorem

Suppose that you have a sum of money to invest in fixed-income securities, and you will not draw from these funds for  $n$  periods (say,  $n$  years). You will invest only in Treasury instruments, and there is a current known spot rate curve for these securities. You have a multitude of choices for structuring a portfolio using your available money. You may select some bonds with long maturities, some zero-coupon bonds, and some bonds with short maturities. If you select a mix of these securities, then, as time passes, you will obtain income from coupons and from the redemption of the short maturity bonds. You may also elect to sell some bonds early, before maturity. As income is generated in these ways, you will reinvest this income in other bonds; again you have a multitude of choices. Finally you will cash out everything at time period  $n$ . How should you invest in order to obtain the maximum amount of money at the terminal time?

To address this question, you must have a model of how interest rates will change in the intervening years, since future rates will determine the prices for bonds that you sell early and those that you buy when reinvesting income. There are a variety of models you could select (some of which might involve randomness, as discussed in Chapter 16), but a straightforward choice is to assume expectations dynamics—so let

us make that assumption. Let us assume that the initial spot rate curve is transformed, after 1 year, to a new curve in accordance with the updating formula presented earlier. This updating is repeated each year. Now, how should you invest?

The answer is revealed by the title of this subsection. It makes absolutely *no* difference how you invest (as long as you remain fully invested). All choices will produce exactly the same result. In particular, investing in a single zero-coupon bond will produce this invariant amount, which is, accordingly,  $(1+s_n)^n$  times your original sum of money. This result is spelled out in the following theorem:

**Invariance theorem** *Suppose that interest rates evolve according to expectations dynamics. Then (assuming a yearly compounding convention) a sum of money invested in the interest rate market for n years will grow by a factor of  $(1+s_n)^n$  independent of the investment and reinvestment strategy (so long as all funds are fully invested).*

**Proof:** The conclusion is easiest to see from the example used earlier. Suppose that  $n = 2$ . You have two basic choices for investment. You can invest in a 2-year zero-coupon bond, or you can invest in a 1-year bond and then reinvest the proceeds at the end of the year. Under expectations dynamics, the reinvestment rate after 1 year will be equal to the current forward rate  $f_{1,2}$ . Both of these choices lead to a growth of  $(1+s_2)^2$ . Any other investment, such as a 2-year bond that makes a coupon payment after 1 year that must be reinvested, will be a combination of these two basic strategies. It should be clear that a similar argument applies for any  $n$ . ■

The simplest way to internalize this result is to think in terms of the short rates. Every investment earns the relevant short rates over its duration. A 10-year zero-coupon bond earns the 10 short rates that are defined initially. An investment rolled over year by year for 10 years earns the 10 short rates that happen to occur. Under expectations dynamics, the short rates do not change; that is, the rate initially implied for a specified period in the future will be realized when that period arrives. Hence no matter how an initial sum is invested, it will progress step by step through each of the short rates.

This theorem is very helpful in discussing how to structure an actual portfolio. It shows that the motivation for selecting a mixture of bonds must be due to anticipated deviations from expectations dynamics—deviations of the realized short rates from their originally implied values. Expectations dynamics is, therefore, in a sense the *simplest* assumption about the future because it implies invariance of portfolio growth with respect to strategy.

## 4.6 Running Present Value

The present value of a cash flow stream is easily calculated in the term structure framework. One simply multiplies each cash flow by the discount factor associated

with the period of the flow and then sums these discounted values; that is, present value is obtained by appropriately discounting all future cash flows.

There is a special, alternative way to arrange the calculations of present value, which is sometimes quite convenient and which has a useful interpretation. This different way is termed **running present value**. It calculates present value in a recursive manner starting with the final cash flow and working backward to the present. This method uses the concepts of expectations dynamics from the previous section, although it is not necessary to assume that interest rates actually follow the expectations dynamics pattern to use the method. Although this method is presented, at this point, as just an alternative to the standard method of calculation, it will be the preferred—indeed standard—method of calculation in later chapters.

To work out the process, suppose  $(x_0, x_1, x_2, \dots, x_n)$  is a cash flow stream. We denote the present value of this stream  $\text{PV}(0)$ , meaning the present value at time zero. Now imagine that  $k$  time periods have passed and we are anticipating the remainder of the cash flow stream, which is  $(x_k, x_{k+1}, \dots, x_n)$ . We could calculate the present value (as viewed at time  $k$ ) using the discount factors that would be applicable then. We denote this present value by  $\text{PV}(k)$ . In general, then, we can imagine the present value running along in time—each period's value being the present value of the remaining stream, but calculated using that period's discount factors. These running values are related to each other in a simple way, which is the basis for the method we describe.

The original present value can be expressed explicitly as

$$\text{PV}(0) = x_0 + d_1 x_1 + d_2 x_2 + \dots + d_n x_n,$$

where the  $d_k$ 's are the discount factors at time zero. This formula can be written in the alternative form

$$\text{PV}(0) = x_0 + d_1 [x_1 + (d_2/d_1)x_2 + \dots + (d_n/d_1)x_n]. \quad (4.2)$$

The values  $d_k/d_1, k = 2, 3, \dots, n$ , are the *inferred discount factors 1 year from now* under an assumption of expectations dynamics (as shown later). Hence,

$$\text{PV}(0) = x_0 + d_1 \text{PV}(1).$$

To show how this works in general, for arbitrary time points, we employ the double-indexing system for discount factors introduced in the previous section. The present values at time  $k$  is

$$\text{PV}(k) = x_k + d_{k,k+1} x_{k+1} + d_{k,k+2} x_{k+2} + \dots + d_{k,n} x_n.$$

Using the discount compounding formula, it follows that  $d_{k,k+j} = d_{k,k+1} d_{k+1,k+j}$ . Hence we may write this equation as

$$\text{PV}(k) = x_k + d_{k,k+1} (x_{k+1} + d_{k+1,k+2} x_{k+2} + \dots + d_{k+1,n} x_n).$$

We can therefore write

$$\text{PV}(k) = x_k + d_{k,k+1} \text{PV}(k+1).$$

This equation states that the present value at time  $k$  is the sum of the current cash flow and a one-period discount of the next present value. Note that  $d_{k,k+1} = 1/(1+f_{k,k+1})$ , where  $f_{k,k+1}$  is the short rate at time  $k$ . Hence in this method discounting always uses short rates to determine the discount factors.

**Present value updating** *The running present values satisfy the recursion*

$$\text{PV}(k) = x_k + d_{k,k+1}\text{PV}(k+1)$$

where  $d_{k,k+1} = 1/(1+f_{k,k+1})$  is the discount factor for the short rate at  $k$ .

To carry out the computation in a recursive manner, the process is initiated by starting at the *final* time. One first calculates  $\text{PV}(n)$  as  $\text{PV}(n) = x_n$  and then  $\text{PV}(n-1) = x_{n-1} + d_{n-1,n}\text{PV}(n)$ , and so forth until  $\text{PV}(0)$  is found.

You can visualize the process in terms of  $n$  people standing strung out, on a time line. You are at the head of the line, at time zero. Each person can observe only the cash flow that occurs at that person's time point. Hence you can observe only the current, time zero, cash flow. How can you compute the present value? Use the running method.

The last person, person  $n$ , computes the present value seen then and passes that value to the first person behind. That person, using the short rate at that time, discounts the value announced by person  $n$ , then adds the observed cash flow at  $n-1$  and passes this new present value back to person  $n-2$ . This process continues, each person discounting according to their short rate, until the running present value is passed to you. Once you hear what the person in front of you announces, you discount it using the initial short rate and add the current cash flow. That is the overall present value.

The running present value  $\text{PV}(k)$  is, of course, somewhat of a fiction. It will be the actual present value of the remaining stream at time  $k$  only if interest rates follow expectations dynamics. Otherwise, entirely different discount rates will apply at that time. However, when computing a present value at time zero, that is, when computing  $\text{PV}(0)$ , the running present value method can be used since it is a mathematical identity.

**Example 4.6 (Constant running rate)** Suppose that the spot rate curve is flat, with  $s_k = r$  for all  $k = 1, 2, \dots, n$ . Let  $(x_0, x_1, x_2, \dots, x_n)$  be a cash flow stream. In the flat case, all inferred forward rates are also equal to  $r$ . (See Exercise 9.) Hence the present value can be calculated as

$$\text{PV}(n) = x_n$$

$$\text{PV}(k) = x_k + \frac{1}{1+r}\text{PV}(k+1).$$

This recursion is run from the terminal time backward to  $k = 0$ .

**Example 4.7 (General running)** A sample present value calculation is shown in Table 4.3. The basic cash flow stream is the first row of the table. We assume that the

**TABLE 4.3**  
**EXAMPLE OF RUNNING PRESENT VALUE**

	Year $k$							
	0	1	2	3	4	5	6	7
Cash flow	20	25	30	35	40	30	20	10
Discount	.943	.935	.93	.926	.923	.921	.917	
PV( $k$ )	168.95	157.96	142.20	120.64	92.49	56.87	29.17	10.00

The present value is found by starting at the final time and working backward, discounting one period at a time.

current term structure is that of Table 4.2, and the appropriate one-period discount rates (found in the first column of the discount factor table in Table 4.2) are listed in the second row of Table 4.3.

The present value at any year  $k$  is computed by multiplying the discount factor listed under that year times the present value of the next year, and then adding the cash flow for year  $k$ . This is done by beginning with the final year and working backward to time zero. Thus we first find  $\text{PV}(7) = 10.00$ . Then  $\text{PV}(6) = 20 + .917 \times 10.00 = 29.17$ ,  $\text{PV}(5) = 30 + .921 \times 29.17 = 56.87$ , and so forth. The present value of the entire stream is  $\text{PV}(0) = 168.95$ .

## 4.7 Floating-Rate Bonds

A floating-rate note or bond has a fixed face value and fixed maturity, but its coupon payments are tied to current (short) rates of interest. Consider, for example, a floating-rate bond that makes coupon payments every 6 months. When the bond is issued, the coupon rate for the first 6 months is set equal to the current 6-month interest rate. At the end of 6 months a coupon payment at that rate is paid; specifically, the coupon is the rate times the face value divided by 2 (because of the 6-month schedule). Then, after that payment, the rate is **reset**: the rate for the next 6 months is set equal to the then current 6-month (short) rate. The process continues until maturity.

Clearly, the exact values of future coupon payments are uncertain until 6 months before they are due. It seems, therefore, that it may be difficult to assess the value of such a bond. In fact at the reset times, the value is easy to deduce—it is equal to par. We highlight this important result.

**Theorem 4.1 (Floating-rate value)** *The value of a floating-rate bond is equal to par at any reset point.*

**Proof:** It is simplest to prove this by working backward using a running present value argument. Look first at the last reset point, 6 months before maturity. We know that the final payment, in 6 months, will be the face value plus the 6-month rate of interest on this amount. The present value at the last reset point is obtained by discounting the total final payment at the 6-month rate—leading to the face value—so the present value is par at that point. Now move back another 6 months

to the previous reset point. The present value there is found by discounting the sum of the next present value and the next coupon payment, again leading to a value of par. We can continue this argument back to time zero. ■

## 4.8 Duration

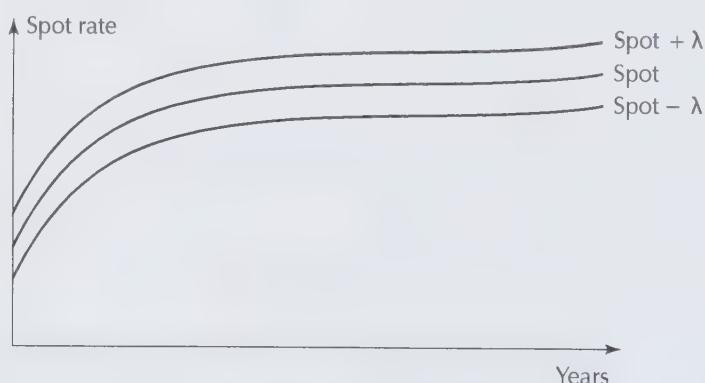
The concept of duration presented in Chapter 3, Section 3.5, can be extended to a term structure framework. We recall that duration is a measure of interest rate sensitivity, which in the earlier development was expressed as sensitivity with respect to yield. In the term structure framework, yield is not a fundamental quantity, but a different, yet similar, measure of risk can be constructed.

The alternative is to consider parallel shifts in the spot rate curve. Specifically, given the spot rates  $s_1, s_2, \dots, s_n$  we imagine that these rates all change together by an additive amount  $\lambda$ . Hence the new spot rates are  $s_1 + \lambda, s_2 + \lambda, \dots, s_n + \lambda$ . This is a hypothetical *instantaneous* change, for the new spot rates are for the same periods as before. This parallel shift of the spot rate curve generalizes a change in the yield because if the spot rate curve were flat, all spot rates would be equal to the common value of yield. Figure 4.3 shows the shifted spot rate curve in the case of a continuous spot rate curve.

Given this notion of a potential change in spot rates, we then can measure the sensitivity of price with respect to the change.

### Fisher-Weil Duration

The details work out most nicely for the case of continuous compounding, and we shall present that case first. Given a cash flow sequence  $(x_{t_0}, x_{t_1}, x_{t_2}, \dots, x_{t_n})$  and the



**FIGURE 4.3 Shifted spot rate curves.** The original spot rate curve is the middle curve. This curve is shifted upward and downward by an amount  $\lambda$  to obtain the other curves. It is possible to immunize a portfolio against such shifts for small values of  $\lambda$ .

spot rate curve  $s_t, t_0 \leq t \leq t_n$ , the present value is

$$PV = \sum_{i=0}^n x_{t_i} e^{-s_{t_i} t_i}.$$

The **Fisher–Weil duration** is then defined as

$$D_{FW} = \frac{1}{PV} \sum_{i=0}^n t_i x_{t_i} e^{-s_{t_i} t_i}.$$

Note that this corresponds exactly to the general definition of duration as a present-value-weighted average of the cash flow times. Clearly  $D_{FW}$  has the units of time and satisfies  $t_0 \leq D \leq t_n$  when all  $x_{t_i} \geq 0$ .

We now consider the sensitivity of price (present value) to a parallel shift of the yield curve and show that it is determined by the Fisher–Weil duration. For arbitrary  $\lambda$  the price is

$$P(\lambda) = \sum_{i=0}^n x_{t_i} e^{-(s_{t_i} + \lambda) t_i}.$$

We then differentiate to find

$$\left. \frac{dP(\lambda)}{d\lambda} \right|_{\lambda=0} = - \sum_{i=0}^n t_i x_{t_i} e^{-s_{t_i} t_i},$$

so immediately we find that the **relative price sensitivity** at  $\lambda = 0$  is

$$\frac{1}{P(0)} \frac{dP(0)}{d\lambda} = -D_{FW}.$$

This essentially duplicates the formula that holds for yield sensitivity presented in Chapter 3.

**Fisher–Weil formulas** *Under continuous compounding, the Fisher–Weil duration of a cash flow stream  $(x_{t_0}, x_{t_1}, \dots, x_{t_n})$  is*

$$D_{FW} = \frac{1}{PV} \sum_{i=0}^n t_i x_{t_i} e^{-s_{t_i} t_i},$$

*where PV denotes the present value of the stream. If all spot rates change to  $s_{t_i} + \lambda, i = 0, 1, 2, \dots, n$ , the corresponding present value function  $P(\lambda)$  satisfies*

$$\frac{1}{P(0)} \frac{dP(0)}{d\lambda} = -D_{FW}.$$

## Discrete-Time Compounding\*

Now we work out the details under the convention of compounding  $m$  times per year. The spot rate in period  $k$  is  $s_k$  (expressed as a yearly rate). Again, we have a cash flow

stream  $(x_0, x_1, x_2, \dots, x_n)$  (where the indexing is by period). The price is

$$P(\lambda) = \sum_{k=0}^n x_k \left(1 + \frac{s_k + \lambda}{m}\right)^{-k}.$$

We then find that

$$\frac{dP(0)}{d\lambda} \equiv \frac{dP(\lambda)}{d\lambda} \Big|_0 = \sum_{k=1}^n -\left(\frac{k}{m}\right) x_k \left(1 + \frac{s_k}{m}\right)^{-(k+1)}.$$

We can relate this to a duration measure by dividing by  $-P(0)$ . Thus we define

$$D_Q \equiv -\frac{1}{P(0)} \frac{dP(0)}{d\lambda} = \frac{\sum_{k=1}^n (k/m)x_k(1+s_k/m)^{-(k+1)}}{\sum_{k=0}^n x_k(1+s_k/m)^{-k}}. \quad (4.3)$$

We term the quantity  $D_Q$  the **quasi-modified duration**. It does have the units of time; however, it is not exactly an average of the cash flow times because  $(1+s_k/m)^{-(k+1)}$  appears in the numerator instead of  $(1+s_k/m)^{-k}$ , which is the discount factor. There is an extra factor of  $(1+s_k/m)^{-1}$  in each numerator term. In the earlier case, where  $s_k$  was constant for all  $k$ , it was possible to pull this extra term outside the summation sign. That led to modified duration. Here such a step is not possible, since the extra factor depends on  $k$ , so we call this rather cumbersome expression by an equally cumbersome name—the quasi-modified duration. It does give the relative price sensitivity to a parallel shift in the spot rate curve. An example is given in the next section.

**Quasi-modified duration** Under compounding  $m$  times per year, the quasi-modified duration of a cash flow stream  $(x_0, x_1, \dots, x_n)$  is

$$D_Q = \frac{1}{PV} \sum_{k=1}^n \left(\frac{k}{m}\right) x_k \left(1 + \frac{s_k}{m}\right)^{-(k+1)}$$

where PV denotes the present value of the stream. If all spot rates change to  $s_k + \lambda$ ,  $k = 1, 2, \dots, n$ , the corresponding present value function  $P(\lambda)$  satisfies

$$\frac{1}{P(0)} \frac{dP(0)}{d\lambda} = -D_Q.$$

Duration is used extensively by investors and professional bond portfolio managers. It serves as a convenient and accurate proxy for interest rate risk. Frequently an institution specifies a guideline that duration should not exceed a certain level, or sometimes a target duration figure is prescribed.

## 4.9 Immunization

The term structure of interest rates leads directly to a new, more robust method for portfolio immunization. This new method does not depend on selecting bonds with

a common yield, as in Chapter 3; indeed, yield does not even enter the calculations. The process is best explained through an example.

**Example 4.8 (A million dollar obligation)** Suppose that we have a \$1 million obligation payable at the end of 5 years, and we wish to invest enough money today to meet this future obligation. We wish to do this in a way that provides a measure of protection against interest rate risk. To solve this problem, we first determine the current spot rate curve. A hypothetical spot rate curve  $s_k$  is shown as the column labeled spot in Table 4.4.

We use a yearly compounding convention in this example in order to save space in the table. We decide to invest in two bonds described as follows:  $B_1$  is a 12-year 6% bond with price 65.95, and  $B_2$  is a 5-year 10% bond with price 101.66. The prices of these bonds are consistent with the spot rates; and the details of the price calculation are given in Table 4.4. The cash flows are multiplied by the discount factors (column  $d$ ), and the results are listed and summed in columns headed  $PV_1$  and  $PV_2$  for the two bonds.

We decide to immunize against a parallel shift in the spot rate curve. We calculate  $dP/d\lambda$ , denoted by  $-PV'$  in Table 4.4, by multiplying each cash flow by  $t$  and by  $(1 + s_t)^{-(t+1)}$  and then summing these. The quasi-modified duration is then the quotient of these two numbers; that is, it equals  $-(1/P)dP/d\lambda$ . The quasi-modified duration of bond 1 is, accordingly,  $466/65.95 = 7.07$ .

We also find the present value of the obligation to be \$627,903.01 and the corresponding quasi-modified duration is  $5/(1 + s_5) = 4.56$ .

To determine the appropriate portfolio we let  $x_1$  and  $x_2$  denote the number of units of bonds 1 and 2, respectively, in the portfolio (assuming, for simplicity, face

**TABLE 4.4  
WORKSHEET FOR IMMUNIZATION PROBLEM**

Year	Spot	$d$	$B_1$	$PV_1$	$-PV'_1$	$B_2$	$PV_2$	$-PV'_2$
1	7.67	.929	6	5.57	5.18	10	9.29	8.63
2	8.27	.853	6	5.12	9.45	10	8.53	15.76
3	8.81	.776	6	4.66	12.84	10	7.76	21.40
4	9.31	.700	6	4.20	15.38	10	7.00	25.63
5	9.75	.628	6	3.77	17.17	110	69.08	314.73
6	10.16	.560	6	3.36	18.29			
7	10.52	.496	6	2.98	18.87			
8	10.85	.439	6	2.63	18.99			
9	11.15	.386	6	2.32	18.76			
10	11.42	.339	6	2.03	18.26			
11	11.67	.297	6	1.78	17.55			
12	11.89	.260	106	27.53	295.26			
Total				65.95	466.00		101.66	386.15
Duration					7.07			3.80

The present values and durations of two bonds are found as transformations of cash flows.

**TABLE 4.5**  
**IMMUNIZATION RESULTS**

	Lambda		
	0	1%	-1%
Bond 1			
Shares	2,208.00	2,208.00	2,208.00
Price	65.94	51.00	70.84
Value	145,602.14	135,805.94	156,420.00
Bond 2			
Shares	4,744.00	4,744.00	4,744.00
Price	101.65	97.89	105.62
Value	482,248.51	464,392.47	501,042.18
Obligation value	627,903.01	600,063.63	657,306.77
Bonds minus obligation	-\$52.37	\$134.78	\$155.40

The overall portfolio of bonds and obligations is immunized against parallel shifts in the spot rate curve.

values of \$100). We then solve the two equations<sup>3</sup>

$$\begin{aligned} P_1x_1 + P_2x_2 &= PV \\ P_1D_1x_1 + P_2D_2x_2 &= PV \times D \end{aligned}$$

where the  $D$ 's are the quasi-modified durations. This leads to  $x_1 = 2,208.17$  and  $x_2 = 4,744.03$ . We round the solutions to determine the portfolio. The results are shown in the first column of Table 4.5, where it is clear that, to within rounding error, the present value condition is met.

To check the immunization properties of this portfolio we change the spot rate curve by adding 1% to each of the spot rate numbers in the first column of Table 4.4. Using these new spot rates, we can again calculate all present values. Likewise, we subtract 1% from the spot rates and calculate present values. The results are shown in the final two columns of Table 4.5. These results show that the immunization property does hold: the change in net present value is only a second-order effect.

Of course, the portfolio is immunized only against parallel shifts in the spot rate curve. It is easy to develop other immunization procedures, which protect against other kinds of shifts as well. Such procedures are discussed in the exercises.

## 4.10 Summary

If observed yield is plotted as a function of time to maturity for a variety of bonds within a fixed risk class, the result is a scatter of points that can be approximated by a curve—the yield curve. This curve typically rises gradually with increasing maturity,

<sup>3</sup> Alternatively, but equivalently, one could solve the equations  $V_1 + V_2 = PV$  and  $D_1V_1 + D_2V_2 = PV \times D$ . Then let  $x_1 = V_1/P_1$  and  $x_2 = V_2/P_2$ .

reflecting the fact that long maturity bonds typically offer higher yields than short maturity bonds. The shape of the yield curve varies continually, and occasionally it may take on an inverted shaped, where yields decrease as the time to maturity increases.

Fixed-income securities are best understood through the concept of the term structure of interest rates. In this structure there is, at any time, a specified interest rate for every maturity date. This is the rate, expressed on an annual basis, that would apply to a zero-coupon bond of the specified maturity. These underlying interest rates are termed spot rates, and if they are plotted as a function of time to maturity, they determine a spot rate curve, similar in character to the yield curve. However, spot rates are fundamental to the whole interest rate market—unlike yields, which depend on the payout pattern of the particular bonds used to calculate them. Once spot rates are determined, it is straightforward to define discount factors for every time, and the present value of a future cash flow is found by discounting that cash flow by the appropriate discount factor. Likewise, the present value of a cash flow stream is found by summing the present values of the individual flow elements.

A series of forward rates can be inferred from a spot rate curve. The forward rate between future times  $t_1$  and  $t_2$  is the interest rate that would be charged for borrowing money at time  $t_1$  and repaying it at time  $t_2$ , but at terms arranged today. These forward rates are important components of term structure theory.

There are three main explanations of the characteristic upward sloping spot rate curve. The first is expectations theory. It asserts that the current implied forward rates for 1 year ahead—that is, the forward rates from year 1 to future dates—are good estimates of next year's spot rates. If these estimates are higher than today's values, the current spot rate curve must slope upward. The second explanation is liquidity preference theory. It asserts that people prefer short-term maturities to long-term maturities because the interest rate risk is lower with short-term maturities. This preference drives up the prices of short-term maturities. The third explanation is the market segmentation theory. According to this theory, there are separate supply and demand forces in every range of maturities, and prices are determined in each range by these forces. Hence the interest rate within any maturity range is more or less independent of that in other ranges. Overall it is believed that the factors in all three of these explanations play a role in the determination of the observed spot rate curve.

Expectations theory forms the basis of the concept of expectations dynamics, which is a particular model of how spot rates might change with time. According to expectations dynamics, next year's spot rates will be equal to the current implied forward rates for 1 year ahead—the rates between year 1 and future years. In other words, the forward rates for 1 year ahead actually will be realized in 1 year. This prediction can be repeated for the next year, and so on. This means that all future spot rates are determined by the set of current forward rates. Expectations dynamics is only a model, and future rates will most likely deviate from the values it delivers; but it provides a logical simple prediction of future rates. As a special case, if the current spot rate curve is flat—say, at 12%—then according to expectations dynamics, the spot rate curve next year will also be flat at 12%. The invariance theorem states that if spot rates evolve according to expectations dynamics, the interest earned on funds

committed to the interest rate market for several years is independent of how those funds are invested.

Present value can be calculated by the running method, which starts from the final cash flow and works backward toward the first cash flow. At any stage  $k$  of the process, the present value is calculated by discounting the next period's present value using the short rate at time  $k$  that is implied by the term structure. This backward moving method of evaluation is fundamental to advanced methods of calculation in various areas of investment science.

Duration can be extended to the term structure framework. The key idea is to consider parallel shifts of the spot rate curve, shifts defined by adding a constant  $\lambda$  to every spot rate. Duration is then defined as  $(-1/P)dP/d\lambda$  evaluated at  $\lambda = 0$ . Fisher-Weil duration is based on continuous-time compounding, which leads to a simple formula. In discrete time, the appropriate, somewhat complicated formula is termed quasi-modified duration.

Once duration is defined, it is possible to extend the process of immunization to the term structure framework. A portfolio of assets designed to fund a stream of obligations can be immunized against a parallel shift in the spot rate curve by matching both the present values and the durations of the assets and the obligations.

## Exercises

- (One forward rate)** If the spot rates for 1 and 2 years are  $s_1 = 6.3\%$  and  $s_2 = 6.9\%$ , what is the forward rate  $f_{1,2}$ ?
- (Spot update)** Given the (yearly) spot rate curve  $s = (5.0, 5.3, 5.6, 5.8, 6.0, 6.1)$ , find the spot rate curve for next year.
- (Construction of a zero)** Consider two 5-year bonds: one has a 9% coupon and sells for 101.00; the other has a 7% coupon and sells for 93.20. Find the price of a 5-year zero-coupon bond.
- (Spot rate project  $\oplus$ )** It is November 5 in the year 2012. The bond quotations of Table 4.6 are available. Assume that all bonds make semiannual coupon payments on the 15th of the month. Estimate the (continuous-time) term structure in the form of a 4th-order polynomial,

$$r(t) = a_0 + a_1 t + a_2 t^2 + a_3 t^3 + a_4 t^4,$$

where  $t$  is time in units of years from today. The discount rate for cash flows at time  $t$  is accordingly  $d(t) = e^{-r(t)t}$ . Recall that accrued interest must be added to the price quoted to get the total price. Estimate the coefficients of the polynomial by minimizing the sum of squared errors between the total price and the price predicted by the estimated term structure curve. Plot the curve and give the five polynomial coefficients.

- (Instantaneous rates  $\diamond$ )** Let  $s(t)$ ,  $0 \leq t \leq \infty$ , denote a spot rate curve; that is, the present value of a dollar to be received at time  $t$  is  $e^{-s(t)t}$ . For  $t_1 < t_2$ , let  $f(t_1, t_2)$  be the forward rate between  $t_1$  and  $t_2$  implied by the given spot rate curve.
  - Find an expression for  $f(t_1, t_2)$ .
  - Let  $r(t) = \lim_{t_2 \rightarrow t} f(t, t_2)$ . We can call  $r(t)$  the instantaneous interest rate at time  $t$ . Show that  $r(t) = s(t) + s'(t)t$ .

**TABLE 4.6**  
**BOND QUOTES**

Coupon	Maturity	Ask price
6.625	15-Feb-22	100.00
9.125	15-Feb-22	100.67
7.875	15-Aug-22	100.69
8.250	15-Aug-22	100.03
8.250	15-Feb-23	100.22
8.375	15-Feb-23	100.38
8.000	15-Aug-23	100.81
8.750	15-Aug-23	102.03
6.875	14-Feb-24	98.16
8.875	14-Feb-24	102.28
6.875	15-Aug-24	97.41
8.625	15-Aug-24	101.72
7.750	15-Feb-25	99.16
11.250	15-Feb-25	109.13
8.500	15-Aug-25	101.41
10.500	15-Aug-25	107.84
7.875	15-Feb-26	99.41
8.875	15-Feb-26	103.00

- (c) Suppose an amount  $x_0$  is invested in a bank account at  $t = 0$  which pays the instantaneous rate of interest  $r(t)$  at all  $t$  (compounded). Then the bank balance  $x(t)$  will satisfy  $dx(t)/dt = r(t)x(t)$ . Find an expression for  $x(t)$ . [Hint: Recall in general that  $ydz + zdy = d(yz)$ .]
6. (Discount conversion) At time zero the one-period discount rates  $d_{0,1}, d_{1,2}, d_{2,3}, \dots, d_{5,6}$  are known to be 0.950, 0.940, 0.932, 0.925, 0.919, 0.913. Find the time zero discount factors  $d_{0,1}, d_{0,2}, \dots, d_{0,6}$ .
7. (Bond taxes) An investor is considering the purchase of 10-year U.S. Treasury bonds and plans to hold them to maturity. Federal taxes on coupons must be paid during the year they are received, and tax must also be paid on the capital gain realized at maturity (defined as the difference between face value and original price). Federal bonds are exempt from state taxes. This investor's federal tax bracket rate is  $t = 30\%$ , as it is for most individuals. There are two bonds that meet the investor's requirements. Bond 1 is a 10-year, 10% bond with a price (in decimal form) of  $P_1 = 92.21$ . Bond 2 is a 10-year, 7% bond with a price of  $P_2 = 75.84$ . Based on the price information contained in those two bonds, the investor would like to compute the theoretical price of a hypothetical 10-year zero-coupon bond that had no coupon payments and required tax payment only at maturity equal in amount to 30% of the realized capital gain (the face value minus the original price). This theoretical price should be such that the price of this bond and those of bonds 1 and 2 are mutually consistent on an after-tax basis. Find this theoretical price, and show that it does not depend on the tax rate  $t$ . (Assume all cash flows occur at the end of each year.)
8. (Real zeros) Actual zero-coupon bonds are taxed as if implied coupon payments were made each year (or really every 6 months), so tax payments are made each year, even though no coupon payments are received. The implied coupon rate for a bond with  $n$  years

to maturity is  $(100 - P_0)/n$ , where  $P_0$  is the purchase price. If the bond is held to maturity, there is no realized capital gain, since all gains are accounted for in the implied coupon payments. Compute the theoretical price of a real 10-year zero-coupon bond. This price is to be consistent on an after-tax basis with the prices of bonds 1 and 2 of Exercise 7.

9. (Flat forwards) Show explicitly that if the spot rate curve is flat [with  $s(k) = r$  for all  $k$ ], then all forward rates also equal  $r$ .
10. (Orange County blues) Orange County managed an investment pool into which several municipalities made short-term investments. A total of \$7.5 billion was invested in this pool, and this money was used to purchase securities. Using these securities as collateral, the pool borrowed \$12.5 billion from Wall Street brokerages, and these funds were used to purchase additional securities. The \$20 billion total was invested primarily in long-term fixed-income securities to obtain a higher yield than the short-term alternatives. Furthermore, as interest rates slowly declined, as they did in 1992–1994, an even greater return was obtained. Things fell apart in 1994, when interest rates rose sharply.

Hypothetically, assume that initially the duration of the invested portfolio was 10 years, the short-term rate was 6%, the average coupon interest on the portfolio was 8.5% of face value, the cost of Wall Street money was 7%, and short-term interest rates were falling at  $\frac{1}{2}\%$  per year.

- (a) What was the rate of return that pool investors obtained during this early period? Does it compare favorably with the 6% that these investors would have obtained by investing normally in short-term securities?
- (b) When interest rates had fallen two percentage points and began increasing at 2% per year, what rate of return was obtained by the pool?
11. (Running PV example) A (yearly) cash flow stream is  $\mathbf{x} = (-40, 10, 10, 10, 10, 10, 10, 10)$ . The spot rates are those of Exercise 2.
  - (a) Find the current discount factors  $d_{0,k}$  and use them to determine the (net) present value of the stream.
  - (b) Find the series of expectations dynamics short-rate discount factors, and use the running present value method to evaluate the stream.
12. (Term structure and bond price) You are given an incomplete specification of the term structure, as specified by the spot rates and forward rates noted next. You also know that the price of a 6-year bond with coupon rate 10% is \$145.749 and the price of a 6-year bond with coupon rate 5% is \$100.315. For all bonds, the face value is \$100, and the coupons are paid annually. Assuming continuous compounding, find the missing rates.

$$s_1 = ?, s_2 = 6.9\%, s_3 = 7.5\%, s_4 = ?, s_5 = 8.4\%, s_6 = ?$$

$$f_{1,2} = 7.8\%, f_{2,3} = 8.7\%, f_{5,6} = ?, f_{1,3} = 8.25\%, f_{2,4} = 11.55\%$$

13. (Duration estimate) A certain bond portfolio has a value of \$1,000 today at a yield of 10%. Yesterday the same portfolio had a value of \$990 at a yield of 10.5%.
  - (a) Estimate what the modified duration was yesterday.
  - (b) Estimate what the Macaulay duration was yesterday, assuming daily compounding (365 days/year).
14. (Pure duration  $\diamond$ ) It is sometimes useful to introduce variations of the spot rates that are different from an additive variation. Let  $\mathbf{s}^0 = (s_1^0, s_2^0, s_3^0, \dots, s_n^0)$  be an initial spot

rate sequence (based on  $m$  periods per year). Let  $s(\lambda) = (s_1, s_2, \dots, s_n)$  be spot rates parameterized by  $\lambda$ , where

$$1 + s_k/m = e^{\lambda/m}(1 + s_k^0/m)$$

for  $k = 1, 2, \dots, n$ . Suppose a bond price  $P(\lambda)$ , is determined by these spot rates. Show that

$$-\frac{1}{P} \frac{dP}{d\lambda} = D$$

is a pure duration; that is, find  $D$  and describe it in words.

- 15.** (Stream immunization  $\oplus$ ) A company faces a stream of obligations over the next 8 years as shown: where the numbers denote thousands of dollars. The spot rate curve is that of Example 4.8. Try to find a portfolio, consisting of the two bonds described in that example, that has the same present value as the obligation stream and is immunized against an additive shift in the spot rate curve.

Year	1	2	3	4	5	6	7	8
	500	900	600	500	100	100	100	50

- 16.** (Mortgage division) Often a mortgage payment stream is divided into a principal payment stream and an interest payment stream, and the two streams are sold separately. We shall examine the component values. Consider a standard mortgage of initial value  $M = M(0)$  with equal periodic payments of amount  $B$ . If the interest rate used is  $r$  per period, then the mortgage principal after the  $k$ th payment satisfies

$$M(k) = (1 + r)M(k - 1) - B$$

for  $k = 0, 1, \dots$ . This equation has the solution

$$M(k) = (1 + r)^k M - \left[ \frac{(1 + r)^k - 1}{r} \right] B.$$

Let us suppose that the mortgage has  $n$  periods and  $B$  is chosen so that  $M(n) = 0$ ; namely,

$$B = \frac{r(1 + r)^n M}{(1 + r)^n - 1}.$$

The  $k$ th payment has an interest component of

$$I(k) = rM(k - 1)$$

and a principal component of

$$P(k) = B - rM(k - 1).$$

- (a) Find the present value  $V$  (at rate  $r$ ) of the principal payment stream in terms of  $B, r, n, M$ .
- (b) Find  $V$  in terms of  $r, n, M$  only.
- (c) What is the present value  $W$  of the interest payment stream?
- (d) What is the value of  $V$  as  $n \rightarrow \infty$ ?
- (e) Which stream do you think has the larger duration—principal or interest?

- 17.** (Short rate sensitivity) Gavin Jones sometimes has flashes of brilliance. He asked his instructor if duration would measure the sensitivity of price to a parallel shift in the short rate curve. (That is,  $r_k \rightarrow r_k + \lambda$ .) His instructor smiled and told him to work it out. He was unsuccessful at first because his formulas became very complicated. Finally he discovered a simple solution based on the running present value method. Specifically, letting  $P_k$  be the present value as seen at time  $k$  and  $S_k = dP_k/d\lambda|_{\lambda=0}$ , the  $S_k$ 's can be found recursively by an equation of the form  $S_{k-1} = -a_k P_k + b_k S_k$ , while the  $P_k$ 's are found by the running method. Find  $a_k$  and  $b_k$ .

## References

For general discussions of term structure theory, see [1–3]. Critical analyses of the expectations explanation are contained in [4] and [5]. The liquidity preference explanation is explored in [6]. Immunization in a term structure environment was originated in [7].

1. Fabozzi, F. J., and F. Modigliani (2008), *Capital Markets: Institutions and Instruments*, Prentice Hall, Englewood Cliffs, NJ.
2. Homer, S., and M. Liebowitz (1972), *Inside the Yield Book: New Tools for Bond Market Strategy*, 4th ed., Prentice Hall, Englewood Cliffs, NJ.
3. Van Home, J. C. (2000), *Financial Market Rates & Flows*, 6th ed., Prentice Hall, Englewood Cliffs, NJ.
4. Russell, S. (July/August 1992), “Understanding the Term Structure of Interest Rates: The Expectations Theory,” *Federal Reserve Bank of St. Louis Review*, 36–51.
5. Cox, J., J. Ingersoll, and S. Ross (September 1981), “A Reexamination of Traditional Hypotheses about the Term Structure of Interest Rates,” *Journal of Finance*, **36**, 769–799.
6. Fama, E. (1984), “The Information in the Term Structure,” *Journal of Financial Economics*, **13**, 509–528.
7. Fisher, L., and R. L. Weil (October 1971), “Coping with the Risk of Market-Rate Fluctuations: Returns to Bondholders from Naive and Optimal Strategies,” *Journal of Business*, **44**, 408–431.

## APPLIED INTEREST RATE ANALYSIS

**U**ltimately, the practical purpose of investment science is to improve the investment process. This process includes identification, selection, combination, and ongoing management. In the ideal case, these process components are integrated and handled as a craft—a craft rooted in scientific principles and meaningful experience, and executed through a combination of intuition and formal problem-solving procedures. This chapter highlights the formal procedures for structuring investments.

The previous chapters provide the groundwork for the analysis of a surprisingly broad set of investment problems. Indeed, interest rate theory alone provides the basis of the vast majority of actual investment studies. Therefore mastery of the previous chapters is adequate preparation to address a wide assortment of investment situations—and appropriate analyses can be conducted with simple practical tools, such as spreadsheet programs, or more complex tools, such as parallel processor computers. To illustrate the range of problems that can be meaningfully treated by the theory developed in earlier chapters, this chapter considers a few typical problem areas. Our treatment of these subjects is only introductory, for indeed there are textbooks devoted to each of these topics. Nevertheless, the solid grounding of the previous chapters allows us to enter these problems at a relatively high level, and to convey quickly the essence of the subject. We consider capital budgeting, bond portfolio construction, management of dynamic investments, and valuation of firms from accounting data. These subjects all represent important investment issues.

To resolve an investment issue with quantitative methods, the issue must first be formulated as a specific problem. There are usually a number of ways to do this, but frequently the best formulation is a version of **optimization**. It is entirely consistent with general investment objectives to try to devise the “ideal” portfolio, to select the “best” combination of projects, to manage an investment to attain the “most favorable” outcome, or to hedge assets to attain the “least” exposure to risk. All of these are, at least loosely, statements of optimization. Indeed, optimization and investment seem like perfect partners. We begin to explore the possibilities of this happy relationship in this chapter.

## 5.1 Capital Budgeting

The capital allocation problem consists of allocating a (usually fixed) budget among a number of investments or projects. We distinguish between **capital budgeting** treated here and **portfolio problems** treated in the next section, although the two are related. Capital budgeting typically refers to allocation among projects or investments for which there are not well-established markets and where the projects are *lumpy* in that they each require discrete lumps of cash (as opposed to securities, where virtually any number of shares can be purchased).

Capital budgeting problems often arise in a firm where several proposed projects compete for funding. The projects may differ considerably in their scale, their cash requirements, and their benefits. The critical point, however, is that even if all proposed projects offer attractive benefits, they cannot all be funded because of a budget limitation. Our earlier study of investment choice, in Chapter 2, focused on situations where the budget was not fixed, and the choice options were mutually exclusive, such as the choice between a red and a green car. In capital budgeting the alternatives may or may not be mutually exclusive, and budget is a definite limitation.

### Independent Projects

The simplest, and classic, type of a capital budgeting problem is that of selecting from a list of independent projects. The projects are independent in the sense that it is reasonable to select any combination from the list. It is not a question of selecting between a red and a green car; we can choose both if we have the required budget. Likewise, the value of one project does not depend on another project also being funded. This standard capital budgeting problem is quite easy to formulate.

Suppose that there are  $m$  potential projects. Let  $b_i$  be the total benefit (usually the net present value) of the  $i$ th project, and let  $c_i$  denote its initial cost. Finally, let  $C$  be the total capital available—the budget. For each  $i = 1, 2, \dots, m$  we introduce the **zero-one variable**  $x_i$ , which is zero if the project is rejected and one if it is accepted.

The problem is then that of solving

$$\begin{aligned} & \text{maximize} \sum_{i=1}^m b_i x_i \\ & \text{subject to} \sum_{i=1}^m c_i x_i \leq C \\ & x_i = 0 \text{ or } 1 \quad \text{for } i = 1, 2, \dots, m. \end{aligned}$$

This is termed a **zero-one programming problem**, since the variables are zero-one variables. It is a formal representation of the fact that projects can either be selected or not, but for those that are selected, both the benefits and the costs are directly additive.

There is an easy way to obtain an approximate solution to this problem, which is quite accurate in many cases. We shall describe this method under the assumption (which can be weakened) that each project requires an initial outlay of funds (a negative cash flow) that is followed by a stream of benefits (a stream of positive cash flows). We define the **benefit-cost ratio** as the ratio of the present worth of the benefits to the magnitude of the initial cost. We then rank projects in terms of this benefit-cost ratio. Projects with the highest ratios offer the best return per dollar invested—the biggest “bang for the buck”—and hence are excellent candidates for inclusion in the final list of selected projects. Once the projects are ranked this way, they are selected one at a time, by order of the ranking, until no additional project can be included without violating the given budget. This method will produce the best value for the amount spent. However, despite this property, the solution found by this approximate method is not always optimal since it may not use the entire available budget. Better solutions may be found by skipping over some high-cost projects so that other projects, with almost as high a benefit-cost ratio, can be included. To obtain true optimality, the zero-one optimization problem can be solved exactly by readily available software programs. However, the simpler method based on the benefit-cost ratio is helpful in a preliminary study. (Some spreadsheet packages have integer programming routines suitable for modest-sized problems.)

**Example 5.1 (A selection problem)** During its annual budget planning meeting, a small computer company has identified several proposals for independent projects that could be initiated in the forthcoming year. These projects include the purchase of equipment, the design of new products, the lease of new facilities, and so forth. The projects all require an initial capital outlay in the coming year. The company management believes that it can make available up to \$500,000 for these projects. The financial aspects of the projects are shown in Table 5.1.

For each project the required initial outlay, the present worth of the benefits (the present value of the remainder of the stream after the initial outlay), and the ratio of these two are shown. The projects are already listed in order of decreasing benefit-cost ratio. According to the approximate method the company would select projects 1, 2, 3, 4, and 5 for a total expenditure of \$370,000 and a total net present value of \$910,000 – \$370,000 = \$540,000. However, this solution is not optimal.

**TABLE 5.1**  
**PROJECT CHOICES**

Project	Outlay (\$1,000)	Present worth (\$1,000)	Benefit-cost ratio
1	100	300	3.00
2	20	50	2.50
3	150	350	2.33
4	50	110	2.20
5	50	100	2.00
6	150	250	1.67
7	150	200	1.33

The outlays are made immediately, and the present worth is the present value of the future benefits. Projects with a high benefit-cost ratio are desirable.

**FIGURE 5.1 Spreadsheet for project choices.** The  $x$ -values are listed in one column. These values are multiplied by the corresponding elements of outlay and net present value to obtain the components of cost and optimal present value in the total package of projects. A zero-one program (within the spreadsheet) adjusts these  $x$ -values to find the optimal set.

Project	Outlay	Present worth	Net PV	Optimal		Optimal PV
				x-value	Cost	
1	100	300	200	1	100	200
2	20	50	30	0	0	0
3	150	350	200	1	150	200
4	50	110	60	1	50	60
5	50	100	50	1	50	50
6	150	250	100	1	150	100
7	150	200	50	0	0	0
<b>Totals</b>					<b>500</b>	<b>610</b>

The proper method of solution is to formulate the problem as a zero-one optimization problem. Accordingly, we define the variables  $x_i$ ,  $i = 1, 2, \dots, 7$ , with  $x_i$  equal to 1 if it is to be selected and 0 if not. The problem is then

$$\begin{aligned} & \text{maximize } 200x_1 + 30x_2 + 200x_3 + 60x_4 + 50x_5 + 100x_6 + 50x_7 \\ & \text{subject to } 100x_1 + 20x_2 + 150x_3 + 50x_4 + 50x_5 + 150x_6 + 150x_7 \leq 500 \\ & \quad x_i = 0 \text{ or } 1 \quad \text{for each } i. \end{aligned}$$

Note that the terms of the objective for maximization are present worth minus outlay—present value.

The problem and its solution are displayed in spreadsheet form in Figure 5.1. It is seen that the solution is to select projects 1, 3, 4, 5, and 6 for a total expenditure of \$500,000 and a total net present value of \$610,000. The approximate method does not account for the fact that using project 2 precludes the use of the more costly, but more beneficial, project 6. Specifically, by replacing 2 by 6 the full budget can be used and, hence, a greater total benefit achieved.

## Interdependent Projects\*

Sometimes various projects are interdependent, the feasibility of one being dependent on whether others are undertaken. We formulate a problem of this type by assuming that there are several independent goals, but each goal has more than one possible method of implementation. It is these implementation alternatives that define the projects. This formulation generalizes the problems studied in Chapter 2, where there was only one goal (such as buying a new car) but several ways to achieve that goal. The more general problem can be treated as a zero-one programming problem.

As an example of the formulation using goals and projects, suppose a transportation authority wishes to construct a road between two cities. Corresponding projects might detail whether the road were concrete or asphalt, two lanes or four, and so forth. Another, independent, goal might be the improvement of a bridge.

In general, assume that there are  $m$  goals and that associated with the  $i$ th goal there are  $n_i$  possible projects. Only one project can be selected for any goal. As before, there is a fixed available budget.

We formulate this problem by introducing the zero-one variables  $x_{ij}$  for  $i = 1, 2, \dots, m$  and  $j = 1, 2, \dots, n_i$ . The variable  $x_{ij}$  equals 1 if goal  $i$  is chosen and implemented by project  $j$ ; otherwise it is 0. The problem is then

$$\begin{aligned} & \text{maximize} \quad \sum_{i=1}^m \sum_{j=1}^{n_i} b_{ij}x_{ij} \\ & \text{subject to} \quad \sum_{i=1}^m \sum_{j=1}^{n_i} c_{ij}x_{ij} \leq C \\ & \quad \sum_{j=1}^{n_i} x_{ij} \leq 1, \quad \text{for } i = 1, 2, \dots, m \\ & \quad x_{ij} = 0 \text{ or } 1 \quad \text{for all } i \text{ and } j. \end{aligned}$$

The exclusivity of the individual projects is captured by the second set of constraints—one constraint for each objective. This constraint states that the sum of the  $x_{ij}$  variables over  $j$  (the sum of the variables corresponding to projects associated with objective  $i$ ) must not exceed 1. Since the variables are all either 0 or 1, this means that at most one  $x_{ij}$  variable can be 1 for any  $i$ . In other words, at most one project associated with goal  $i$  can be chosen.

In general this is a more difficult zero-one programming problem than that for independent projects. This new problem has more constraints, hence it is not easy to obtain a solution by inspection. In particular, the approximate solution based on benefit-cost ratios is not applicable. However, even large-scale problems of this type can be readily solved with modern computers.

**Example 5.2 (County transportation choices)** Suppose that the goals and specific projects shown in Table 5.2 are being considered by the County Transportation Authority.

**TABLE 5.2**  
**TRANSPORTATION ALTERNATIVES**

	Cost (\$1,000)	NPV (\$1,000)
Road between Augen and Burger		
1 Concrete, 2 lanes	2,000	4,000
2 Concrete, 4 lanes	3,000	5,000
3 Asphalt, 2 lanes	1,500	3,000
4 Asphalt, 4 lanes	2,200	4,300
Bridge at Cay Road		
5 Repair existing	500	1,000
6 Add lane	1,500	1,500
7 New structure	2,500	2,500
Traffic Control in Downsberg		
8 Traffic lights	100	300
9 Turn lanes	600	1,000
10 Underpass	1,000	2,000

*At most one project can be selected for each major objective.*

There are three independent goals and a total of 10 projects. Table 5.2 shows the cost and the net present value (after the cost has been deducted) for each of the projects. The total available budget is \$5 million. To formulate this problem we introduce a zero–one variable for each project. (However, for simplicity we index these variables consecutively from 1 through 10, rather than using the double indexing procedure of the general formulation presented earlier.) The problem formulation can be expressed as

$$\begin{aligned}
 & \text{maximize} && 4x_1 + 5x_2 + 3x_3 + 4.3x_4 + x_5 + 1.5x_6 + 2.5x_7 + .3x_8 + x_9 + 2x_{10} \\
 & \text{subject to} && 2x_1 + 3x_2 + 1.5x_3 + 2.2x_4 + .5x_5 + 1.5x_6 + 2.5x_7 + .1x_8 + .6x_9 + x_{10} \leq 5 \\
 & && x_1 + x_2 + x_3 + x_4 \leq 1 \\
 & && x_5 + x_6 + x_7 \leq 1 \\
 & && x_8 + x_9 + x_{10} \leq 1 \\
 & && x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10} = 0 \text{ or } 1.
 \end{aligned}$$

This problem and its solution are clearly displayed by a spreadsheet, as illustrated in Figure 5.2. The solution is that projects 2, 5, and 10 should be selected, for a cost of \$4,500,000 and a total present value of \$8,000,000.

This method for treating dependencies among projects can be extended to situations where precedence relations apply (that is, where one project cannot be chosen unless another is also chosen) and to capital budgeting problems with additional financial constraints. Typically these more general problems merely impose additional constraints among the variables.

Project	Cost	NPV	Optimal	Cost	NPV	Goals
	(\$1,000)	(\$1,000)	x-values			
1 Concrete, 2 lanes	2,000	4,000	0	0	0	
2 Concrete, 4 lanes	3,000	5,000	1	3,000	5,000	
3 Asphalt, 2 lanes	1,500	3,000	0	0	0	
4 Asphalt, 4 lanes	2,200	4,300	0	0		1
5 Repair existing	500	1,000	1	500	1,000	
6 Add lane	1,500	1,500	0	0	0	
7 New structure	2,500	2,000	0	0	0	1
8 Traffic lights	100	300	0	0	0	
9 Turn lanes	600	1,000	0	0	0	
10 Underpass	1,000	2,000	1	1,000	2,000	1
Totals				4,500	8,000	

**FIGURE 5.2 Transportation spreadsheet.** The  $x$ -values are shown in one column; the corresponding elements of cost and net present value in the next columns. Also, the number of projects included for each goal are shown in the final column. These numbers are constrained to be less than or equal to 1. The optimal  $x$ -values are found by a zero-one programming package.

Although capital budgeting is a useful concept, its basic formulation is somewhat flawed. The *hard* budget constraint is inconsistent with the underlying assumption that it is possible for the investor (or organization) to borrow unlimited funds at a given interest rate. Indeed, in theory one should carry out *all* projects that have positive net present value. In practice, however, the assumption that an unlimited supply of capital is available at a fixed interest rate does not hold. A bank may impose a limited credit line, or in a large organization investment decisions may be decentralized by passing down budgets to individual organizational units. It is therefore often useful to in fact solve the capital budgeting problem. However, it is usually worth solving the problem for various values of the budget to measure the sensitivity of the benefit to the budget level.

## 5.2 Optimal Portfolios

Portfolio optimization is another capital allocation problem, similar to capital budgeting. The term **optimal portfolio** usually refers to the construction of a portfolio of financial securities. However, the term is also used more generally to refer to the construction of any portfolio of financial assets, including a “portfolio” of projects. When the assets are freely traded in a market, certain pricing relations apply that may not apply to more general, nontraded assets. This feature is an important distinction that is highlighted by using the term **portfolio optimization** for problems involving securities.

This section considers only portfolios of fixed-income instruments. As we know, a fixed-income instrument that returns cash at known points in time can be described by listing the stream of promised cash payments (and future cash outflows, if any). Such an instrument can be thought of as corresponding to a list or a vector, with the payments as components, defining an associated cash flow stream. A portfolio is just a combination of such streams, and can be represented as a combination of the individual lists or vectors representing the securities. Spreadsheets offer one convenient way to handle such combinations.

## The Cash Matching Problem

A simple optimal portfolio problem is the **cash matching problem**. To describe this problem, suppose that we face a known sequence of future monetary obligations. (If we manage a pension fund, these obligations might represent required annuity payments.) We wish to invest now so that these obligations can be met as they occur; and accordingly, we plan to purchase bonds of various maturities and use the coupon payments and redemption values to meet the obligations. The simplest approach is to design a portfolio that will, without future alteration, provide the necessary cash as required.

To formulate this problem mathematically, we first establish a basic time period length, with cash flows occurring at the end of these periods. For example, we might use 6-month periods. Our obligation is then a stream  $\mathbf{y} = (y_1, y_2, \dots, y_n)$ , starting one period from now. (We use boldface letters to denote an entire stream.) Likewise each bond has an associated cash flow stream of receipts, starting one period from now. If there are  $m$  bonds, we denote the stream associated with one unit of bond  $j$  by  $\mathbf{c}_j = (c_{1j}, c_{2j}, \dots, c_{nj})$ . The price of bond  $j$  is denoted by  $p_j$ . We denote by  $x_j$  the amount of bond  $j$  to be held in the portfolio. The cash matching problem is to find the  $x_j$ 's of minimum total cost that guarantee that the obligations can be met. Specifically,

$$\begin{aligned} & \text{minimize} \sum_{j=1}^m p_j x_j \\ & \text{subject to} \sum_{j=1}^m c_{ij} x_j \geq y_i \quad \text{for } i = 1, 2, \dots, n \\ & \quad x_j \geq 0 \quad \text{for } j = 1, 2, \dots, m. \end{aligned}$$

The objective function to be minimized is the total cost of the portfolio, which is equal to the sum of the prices of the bonds times the amounts purchased. The main set of constraints are the cash matching constraints. For a given  $i$  the corresponding constraint states that the total amount of cash generated in period  $i$  from all  $m$  bonds must be at least equal to the obligation in period  $i$ . The final constraint rules out the possibility of selling bonds short.

**TABLE 5.3**  
**CASH MATCHING EXAMPLE**

Yr	Bonds										Req'd	Actual
	1	2	3	4	5	6	7	8	9	10		
1	10	7	8	6	7	5	10	8	7	100	100	171.74
2	10	7	8	6	7	5	10	8	107		200	200.00
3	10	7	8	6	7	5	110	108			800	800.00
4	10	7	8	6	7	105					100	119.34
5	10	7	8	106	107						800	800.00
6	110	107	108								1,200	1,200.00
<i>p</i>	109	94.8	99.5	93.1	97.2	92.9	110	104	102	95.2	2,381.14	
<i>x</i>	0	11.2	0	6.81	0	0	0	6.3	0.28	0	Cost	

A spreadsheet layout clearly shows the problem and its solution. In this example, the cash flow streams of 10 different bonds are shown, year by year, as 10 columns in the array. The current price of each bond is listed below the stream, and the amount to be included in a portfolio is listed below the price. Cash flows required to be generated by the portfolio are shown in the penultimate column, and those actually generated are shown in the last column.

This problem can be clearly visualized in terms of an array of numbers in a spreadsheet, as in the following example.

**Example 5.3 (A 6-year match)** We wish to match cash obligations over a 6-year period. We select 10 bonds for this purpose (and for simplicity all accounting is done on a yearly basis). The cash flow structure of each bond is shown in the corresponding column in Table 5.3. Below this column is the bond's current price. For example, the first column represents a 10% bond that matures in 6 years. This bond is selling at 109. The second to last column shows the yearly cash requirements (or obligations) for cash to be generated by the portfolio. We formulate the standard cash matching problem as a linear programming problem and solve for the optimal portfolio. (The solution can be found easily by use of a standard linear programming package such as those available on some spreadsheet programs.) The solution is given in the bottom row of Table 5.3. The actual cash generated by the portfolio is shown in the right-hand column. This column is computed by multiplying each bond column  $j$  by its solution value  $x_j$  and then summing these results. The minimum total cost of the portfolio is also indicated in the table.

Note that in two of the years extra cash, beyond what is required, is generated. This is because there are high requirements in some years, and so a large number of bonds must be purchased that mature at those dates. However, these bonds generate coupon payments in earlier years and only a portion of these payments is needed to meet obligations in those early years. A smoother set of cash requirements would not lead to such surpluses.

There is a fundamental flaw in the cash matching problem as formulated here, as evidenced by the surpluses generated in our example. The surpluses amount to

extra cash, which is essentially thrown away since it is not used to meet obligations and is not reinvested. In reality, such surpluses would be immediately reinvested in instruments that were available at that time. Such reinvestment can be accommodated by a slight modification of the problem formulation, but some assumptions about the nature of future investment opportunities must be introduced. The simplest is to assume that extra cash can be carried forward at zero interest; that it can, so to speak, be put under the mattress to be recovered when needed. This flexibility is introduced by adjoining artificial “bonds” having cash flow streams of the form  $(0, \dots, 0, -1, 1, 0, \dots, 0)$ . Such a bond is “purchased” in the year with the  $-1$  (since it absorbs cash) and is “redeemed” the next year. An even better formulation would allow surplus cash to be invested in actual bonds, but to incorporate this feature an assumption about future interest rates (or, equivalently, about future bond prices) must be made. One logical approach is to assume that prices follow expectations dynamics based on the current spot rate curve. Then if  $r'$  is the estimate of what the 1-year interest rate will be a year from now, which under expectations dynamics is the current forward rate  $f_{1,2}$ , a bond of the form  $(0, -1, 1 + r', 0, \dots, 0)$  would be introduced. The addition of such future bonds allows surpluses to be reinvested, and this addition will lead to a different solution than the simple cash matching solution given earlier.

Other modifications to the basic cash matching problem are possible. For example, if the sums involved are not large, then account might be made of the integer nature of the required solution; that is, the  $x_i$  variables might be restricted to be integers. Other modifications combine immunization with cash matching.

## 5.3 Dynamic Cash Flow Processes

To produce excellent results, many investments require deliberate ongoing management. For example, the course of a project within a firm might be guided by a series of operational decisions. Likewise, a portfolio of financial instruments might (and should) be modified systematically over time. The selection of an appropriate sequence of actions that affect an investment’s cash flow stream is the problem of dynamic management.

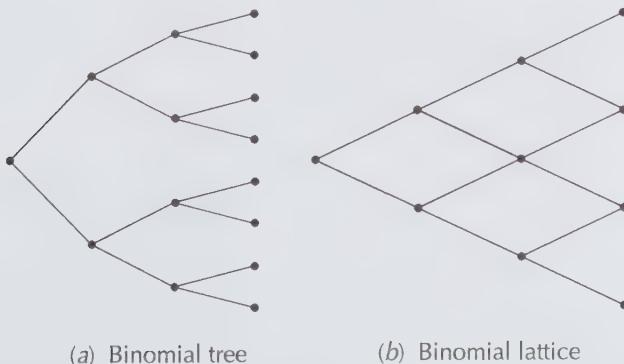
Imagine, for example, that you have purchased an oil well. This is an investment project, and to obtain good results from it, it must be carefully managed. In this case you must decide, each month, whether to pump oil from your well or not. If you do pump oil, you will incur operational costs and receive revenue from the sale of oil, leading to a profit; but you will also reduce the oil reserves. Your current pumping decision clearly influences the future possibilities of production. If you believe that current oil prices are low, you may wisely choose not to pump now, but rather to save the oil for a time of higher prices.

Discussion of this type of problem within the context of deterministic cash flow streams is especially useful—both because it is an important class of problems, and because the method used to solve these problems, **dynamic programming**, is used also in Part 3 of the book. This simpler setting provides a good foundation for that later work.

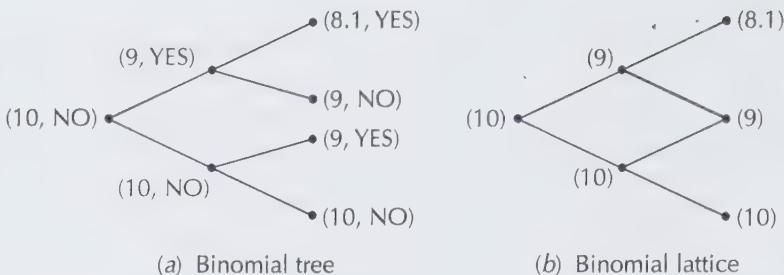
## Representation of Dynamic Choice

A deterministic investment is defined by its cash flow stream, say,  $\mathbf{x} = (x_0, x_1, x_2, \dots, x_n)$ , but the magnitudes of the cash flows in this stream often depend on management choices in a complex fashion. In order to solve dynamic management problems, we need a way to represent the possible choices at each period, and the effect that those choices have on future cash flows. In short, we need a **dynamic model**. There are several mathematical structures that can be used to construct such a model, but the simplest is a **graph**. In this structure, the time points at which cash flows occur are represented by points along the horizontal direction, as usual. In the vertical direction above each such time point is laid out a set of **nodes**, which represent the different possible **states** or conditions of the process at that time. Nodes from one time to the next are connected by **branches** or **arcs**. A branch represents a possible path from a node at one time to another node at the next time. Different branches correspond to different management actions, which guide the course of the process. Simple examples of such graphs are that of a **binomial tree** and a **binomial lattice**, illustrated in Figure 5.3(a) and (b). In such a tree there are exactly two branches leaving each node. The leftmost node corresponds to the situation at the initial time, the next vertical pair of nodes represent the two possibilities at time 1, and so forth. (In the figure only four time points are shown.)

The best way to describe the meaning of the tree is to walk through an example. Let us again consider the management of the oil well you recently purchased. At any time you can either pump oil or not. A node in the tree represents the condition of the well, defined by the size of its reserves, the state of repair, and so forth. To model your choices as a tree, you should start at the leftmost node of the tree, which represents the initial condition of the well. You have only two choices at that time: pump or don't pump. Assign one of these choices to an upward movement and the other to a downward movement; suppose that pumping corresponds to moving upward and nonpumping corresponds to moving downward. At the next time point your well is at one of the two nodes for that time. Again you make a choice and move either



**FIGURE 5.3 Graph representations.** A tree is a general way to represent dynamic choice.



**FIGURE 5.4** Trees showing oil well states. Pumping corresponds to an upward movement; no pumping corresponds to a downward movement. The tree in (a) accounts for both the level of reserves and the status of a crew. If only the reserve levels affect the profit, some nodes combine, forming a binomial lattice, as shown in (b).

up or down. As you make your decisions, you move through the tree, from left to right, from node to node, along a particular path of branches. The path is uniquely determined by your choices; that is, the condition of the well through time and the magnitude of your overall profit are determined by your choices and represented by this unique path through the tree.

Suppose, specifically, that the well has initial reserves of 10 million barrels of oil. Each year it is possible to pump out 10% of the current reserves, but to do so a crew must be hired and paid. However, if a crew is already on hand, because it was used in the previous year, the hiring expenses are avoided. Therefore, to calculate the profit that can be obtained in any year, it is necessary to know the level of oil reserves and whether a crew is already on hand. Hence we label each node of the tree showing the reserve level and the status of a crew. For example, the label (9, YES) means that the reserves are 9,000,000 barrels and there is a crew on hand. A complete tree for the two periods is shown in Figure 5.4(a).

If crews can be assembled with no hiring cost, it is not necessary to keep track of the crew status. We can therefore drop one component from the node labels and keep only the reserve level. If we do that, some nodes that had distinct labels in the original tree will now have identical labels. In the example illustrated in Figure 5.4, two of the nodes at the final time both have a reserve level of 9 (meaning 9 million barrels). Since the labels are identical, we can combine these nodes into a single node, as shown in Figure 5.4(b). If the tree were extended for additional time periods, this combining effect would happen frequently, and as a result the tree could be collapsed to a binomial lattice. A typical binomial lattice is shown in Figure 5.3(b). In such a graph, moving up and then down leads to the same node as moving down and then up. There are fewer nodes in a binomial lattice than in a binomial tree.

In terms of the oil well, if the only relevant factor for determining profit is the reserve level, it is clear that starting at any node, an upward movement in the tree (corresponding to pumping) followed by a downward movement (corresponding to not pumping) is identical in its influence on reserves to a downward movement followed by an upward movement. Both combinations deplete the reserves by the

same amount. Hence a binomial lattice can be used to represent the management choices, as in Figure 5.4(b).

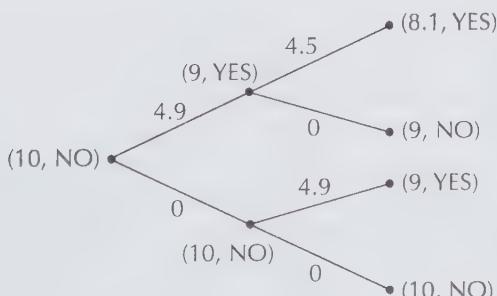
We used a binomial tree or a binomial lattice for the oil well example, which is appropriate when there are only two possible choices at each time. If there were three choices, we could form a **trinomial tree** or a **trinomial lattice**, having three branches emanating from each node. Clearly, any finite number of choices can be accommodated. (It is only reasonable to draw small trees on paper, but a computer can handle larger trees quite effectively, up to a point.)

## Cash Flows in Graphs

The description of the nodes of a graph as states of a process is only an intermediate step in the representation of a dynamic investment situation. The essential part of the final representation is an assignment of cash flows to the various branches of the graph. These cash flows are used to evaluate management alternatives.

In the first oil well example, where crew hiring costs are not zero, suppose that the cost of hiring a crew is \$100,000. (This represents just the initial hiring cost, not the wages paid.) Suppose the profit from oil production is \$5.00 per barrel. Finally, suppose that at the beginning of a year the level of reserves in the well is  $x$ . Then the net profit for a year of production is  $\$5 \times .10 \times x - \$100,000$  if a crew must be hired, and  $\$5 \times .10 \times x$  if a crew is already on hand. We can enter these values on the branches of the tree, indicating that that much profit is attained if that branch is selected. These values are shown in Figure 5.5 in units of millions of dollars.

Since only the cash flow values on the branches are important for analysis, it would be possible (conceptually) to bypass the step of describing the nodes as states of the process. However, in practice the node description is important because the cash flow values are determined from these descriptions by an accounting formula. If someone gave us the tree with cash flow values specified on all branches, that would be sufficient; we would not need the node descriptions. In practice, someone must first characterize the nodes, as we did earlier, so that the cash flows can be determined.



**FIGURE 5.5 Oil well cash flow tree.** The cash flow corresponding to a decision is listed on the branch corresponding to that decision. These cash flow values are determined by the node state and the decision.

In representations of this kind it must also be stated whether the cash flow of a branch occurs at the beginning or at the end of the corresponding time period. In reality, a branch cash flow is often spread out over the entire period, but the model assigns a lump value at one end or the other (or sometimes a part at the beginning and another part at the end). The choice may vary with the situation being represented.

In some cases there is cash flow associated with the termination of the process, whose value varies with the final node achieved. This is a **final reward** or **salvage value**. These values are placed on the graph at the corresponding final nodes. In the oil well example, the final value might be the value for which the well could be sold.

## 5.4 Optimal Management

Once we have a graph representation of the cash flow process associated with an investment, we can apply the principles of earlier chapters to determine the optimal management plan. Each path through the tree determines a specific cash flow stream; hence it is only necessary to select the path that is best. Usually this is the path that has the largest present value. So one way to solve the problem is to list all the possible streams, corresponding to all the possible paths, compute their respective present values, and select the largest one. We then manage the investment by following the path that corresponds to that maximal present value.

Although this method will work well for small problems, it is plagued by the **curse of dimensionality** for large problems. The number of possible paths in a tree grows exponentially with the number of periods. For example, in an  $n$ -period binomial tree (with nodes at time  $0, 1, 2, \dots, n$ ) the number of paths is  $2^n$ . So if  $n = 12$  (say, 1 year of monthly decisions), there are 4,096 possible paths. And if there were 10 possible choices each month, this figure would rise to  $10^{12}$ , which is beyond the capability of straightforward computation. We can use the computational procedure of dynamic programming to search much more efficiently.

## Running Dynamic Programming

Dynamic programming solves a problem step by step, starting at the termination time and working back to the beginning. For this reason, dynamic programming is sometimes characterized by the phrase, “it solves the problem backward.”

A special version of dynamic programming, based on the running present value method of Section 4.6, is especially convenient for investment problems. We call this method **running dynamic programming**. It is the method that we develop here and that is used throughout the text.

Suppose an investment with a dynamic cash flow is represented by a graph as described earlier. For simplicity, we assume periods are 1 year in length, and we use yearly compounding. A path through the graph generates a cash flow stream  $c_0, c_1, \dots, c_{n-1}$  (with each flow occurring at the beginning of the period), corresponding to the arcs that it passes along, and the path also determines a termination flow

$V_n$  at the final node. The present value of this complete stream is

$$PV = c_0 + \frac{c_1}{1+s_1} + \frac{c_2}{(1+s_2)^2} + \cdots + \frac{c_{n-1}}{(1+s_{n-1})^{n-1}} + \frac{V_n}{(1+s_n)^n},$$

where the  $s_k$ 's are the spot rates. A path is defined by a particular series of decisions—one choice at each node. We wish to determine those choices that maximize the resulting present value.

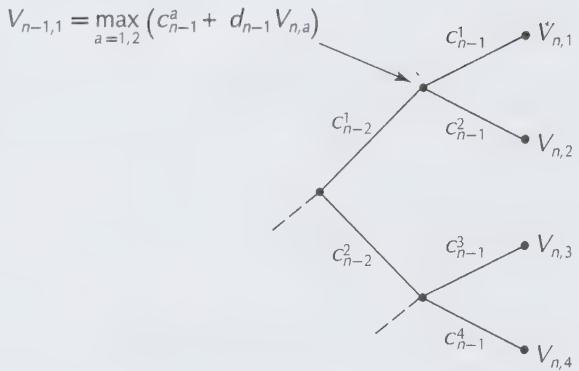
In the running method, we use the one-period discount factors  $d_k = 1/(1+r_k)$ , where  $r_k$  is the short rate  $r_k = f_{k,k+1}$ , and we evaluate the present value step by step in backward order. In particular, in running dynamic programming we assign to each node a value equal to the best running present value that can be obtained from that node, neglecting all previous cash flows. For the  $i$ th node at time  $k$ , denoted by  $(k,i)$ , the best running value is called  $V_{ki}$ . We refer to these values as  $V$ -values.

The  $V$ -values at the final nodes are just the terminal values of the investment process. These values are clearly the present values—as seen at time  $n$ —that can be attained neglecting the past. Hence the  $V$ -values at the final nodes are already given as part of the problem description.

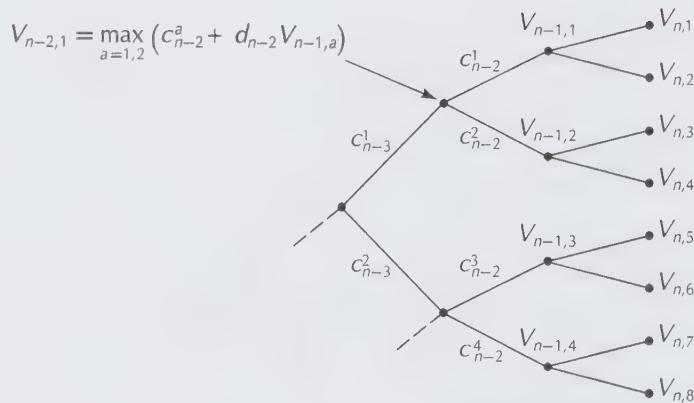
The dynamic programming procedure next addresses the nodes at time  $n-1$ . For any node  $i$  at time  $n-1$ , we pretend that the underlying investment process has taken us to that node. The decisions for previous nodes have already been made, and the corresponding previous cash flows  $c_0, c_1, \dots, c_{n-2}$  have already occurred. Only one decision remains: we must determine which arc to follow from node  $(n-1, i)$  to some final node at time  $n$ . Since we can do nothing about past decisions (in this pretending viewpoint), it is clear that we should select the arc that maximizes the present value as seen at time  $n-1$  (the running present value). Specifically, if we index the arcs by the node number  $a$  they reach at time  $n$ , we should look at the values  $c_{n-1}^a + d_n V_{n,a}$ . (Here  $c_{n-1}^a$  is the cash flow associated with arc  $a$  and  $V_{n,a}$  is the  $V$ -value at the node to which arc  $a$  leads.) After calculating these sums for every arc  $a$  emanating from node  $(n-1, i)$ , we select the largest of these values and denote that value by  $V_{n-1,i}$ . This is the best running present value that can be attained from node  $(n-1, i)$ ; and hence it is the correct  $V$ -value. This procedure, illustrated in Figure 5.6, is repeated for each of the nodes at time  $n-1$ .

Next the same procedure is carried out at time  $n-2$ . We assume that the investment process is at a particular node  $(n-2, i)$ . Each branch  $a$  emanating from that node produces a cash flow and takes the process to a corresponding node  $a$  at time  $n-1$ . If  $c_{n-2}^a$  is the cash flow associated with this choice, the total contribution to (running) present value, accounting for the future as well, is  $c_{n-2}^a + d_{n-1} V_{n-1,a}$  because the running present value is equal to the current cash flow plus a discounted version of the running present value of the next period. We compute these new values for all possible arcs and select the largest. This maximal value is defined to be  $V_{n-2,i}$ . This procedure, illustrated in Figure 5.7, is carried out for every node at time  $n-2$ .

This procedure is continued, working backward until time zero is reached, where there is only one node. The  $V$ -value determined there is the optimal present value as seen at time zero, and hence it is the overall best value. The optimal decisions and cash flows can easily be determined as a by-product of the dynamic programming



**FIGURE 5.6 First recursive step of dynamic programming.** Assuming that the first  $n - 1$  steps of the process have been completed, we evaluate the best that can be done for the last step. For any node at time  $n - 1$  we find the maximum running present value from that node.



**FIGURE 5.7 Second stage of dynamic programming.** Assuming that the first  $n - 2$  stages of the process have been completed, we evaluate the best running present value for the remaining two stages.

procedure, either by recording them at the nodes as the  $V$ -values are computed, or by working forward, using the known future  $V$ -values.

The running dynamic programming method can be written very succinctly by a recurrence relation. Define  $c_{ki}^a$  to be the cash flow generated by moving from node  $(k, i)$  to node  $(k + 1, a)$ . The recursion procedure is

$$V_{ki} = \underset{a}{\text{maximize}}(c_{ki}^a + d_k V_{k+1,a}).$$

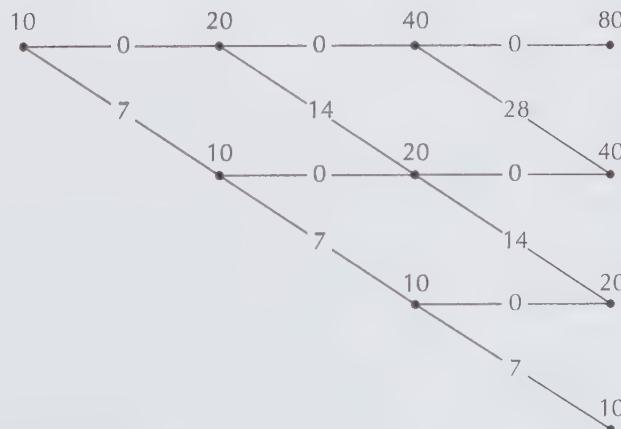
An example will make all of this clear.

## Examples

**Example 5.4 (Fishing problem)** Suppose that you own both a lake and a fishing boat as an investment package. You plan to profit by taking fish from the lake. Each season you decide either to fish or not to fish. If you do not fish, the fish population in the lake will flourish, and in fact it will double by the start of the next season. If you do fish, you will extract 70% of the fish that were in the lake at the beginning of the season. The fish that were not caught (and some before they are caught) will reproduce, and the fish population at the beginning of the next season will be the same as at the beginning of the current season. So corresponding to whether you abstain or fish, the fish population will either double or remain the same, and you get either nothing or 70% of the beginning-season fish population. The initial fish population is 10 tons. Your profit is \$1 per ton. The interest rate is constant at 25%, which means that the discount factor is .8 each year. Unfortunately you have only three seasons to fish. The management problem is that of determining in which of those seasons you should fish.

The situation can be described by the binomial lattice shown in Figure 5.8. The nodes are marked with the fish population. A lattice, rather than a tree, is appropriate because only the fish population in the lake is relevant at any time. The manner by which that population was achieved has no effect on future cash flows. The value on a branch indicates the catch (and hence the cash flow) associated with that branch. Horizontal branches correspond to no fishing and no catch, whereas downward directed branches correspond to fishing.

The problem is solved by working backward. We assign the value of 0 to each of the final nodes, since once we are there we can no longer fish. Then at each of the nodes one step from the end we determine the maximum possible cash flow. (Clearly, we fish in every case.) This determines the cash flow received that season, and we



**FIGURE 5.8 Fishing problem.** The node values are the tonnage of fish in the lake; the branch values are cash flows.

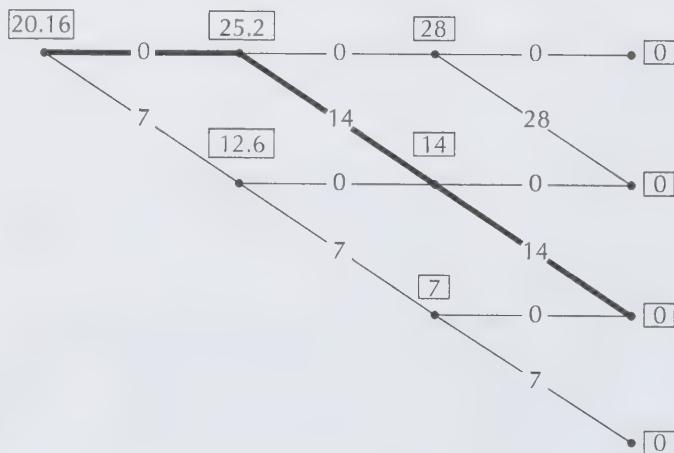
assume that we obtain that cash at the beginning of the season. Hence we do *not* discount the profit. The value obtained is the (running) present value, as viewed from that time. These values are indicated on a copy of the lattice in Figure 5.9.

Next we back up one time period and calculate the maximum present values at that time. For example, for the node just to the right of the initial node, we have

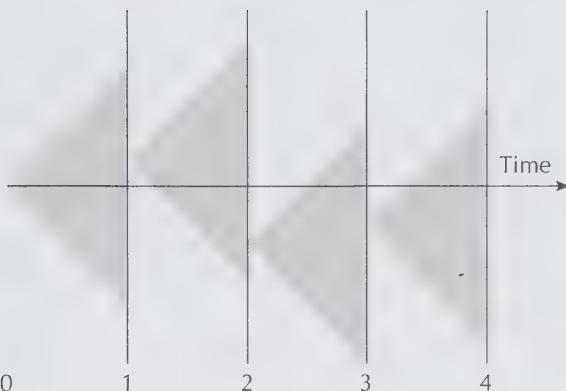
$$V = \max (.8 \times 28, 14 + .8 \times 14).$$

The maximum is attained by the second choice, corresponding to the downward branch, and hence  $V = 14 + .8 \times 14 = 25.2$ . The discount rate of  $1/1.25 = .8$  is applicable at every stage since the spot rate curve is flat. (See Section 4.6.) Finally, a similar calculation is carried out for the initial node. The value there gives the maximum present value. The optimal path is the path determined by the optimal choices we discovered in the procedure. The optimal path for this example is indicated in Figure 5.9 by the heavy line. In words, the solution is not to fish the first season (to let the fish population increase) and then fish the next two seasons (to harvest the population).

The lattice structure can accommodate any finite number of branches emanating from a node. The limit of this kind of construction is a continuous lattice, having a continuum of nodes at any stage and a continuum of possible decisions at any node. For example, in the case of the oil well discussed in the previous section, from a total reserve  $R$  you might pump any amount  $z$  between, say, 0 and  $M$ , leading to a new reserve of  $R - z$ . Your choice  $z$  is continuous, and so is the level of reserves. This type of lattice is illustrated schematically in Figure 5.10. Here each vertical line represents the continuum of nodes possible at a particular time. (At the initial time there is only one node.) The fan emanating from a node represents the fan of



**FIGURE 5.9 Calculations for fish problem.** The node values are now the optimal running present values, found by working backward from the terminal nodes. The branch values are cash flows.



**FIGURE 5.10 Continuous lattice.** A continuous lattice is a powerful way to represent situations where there is a continuum of possible choices every period.

possibilities for traveling to a subsequent node. Only one fan is indicated for each time, whereas actually there is such a fan emanating from every point on the vertical line. This dynamic structure works very much like the finite-node case: The process starts at the initial node, and one of the possible choices is selected. This leads to a specific node point on the line for the next time, and the process continues. Optimizing such a process by dynamic programming works in the reverse direction, just like in the finite case, but is made more difficult by the fact that a  $V$ -value must be assigned to every point on each node line. Hence  $V$  is a function defined on the line. In some cases this function has a simple analytic form, and then the dynamic programming procedure can be carried out explicitly. An illustration of this kind is shown in the next example, which, by the way, is the next in our continuing sequence of gold mine examples.

**Example 5.5 (Complexico mine\*)** The Complexico mine is for lease. This mine has been worked heavily and is approaching depletion. It is becoming increasingly difficult to extract rich ore. In fact, if  $x$  is the amount of gold remaining in the mine at the beginning of a year, the cost to extract  $z < x$  ounces of gold in that year is  $\$500z^2/x$ . (Note that as  $x$  decreases, it becomes more difficult to obtain gold.) It is estimated that the current amount of gold remaining in the mine is  $x_0 = 50,000$  ounces. The price of gold is \$400/oz. We are contemplating the purchase of a 10-year lease of the Complexico mine. The interest rate is 10%. How much is this lease worth?

To solve this problem we must know how to operate the mine optimally over the 10-year period. In particular, we must determine how much gold to mine each year in order to obtain the maximum present value. To find this optimal operating plan, we represent the mine by a continuous lattice, with the nodes at any time representing the amount of gold remaining in the mine at the beginning of that year. We denote this amount by  $x$ . This amount determines the optimal value of the remaining lease from that point on.

We index the time points by the number of years since the beginning of the lease. The initial time is 0, the end of the first year is 1, and so forth. The end of the lease is time 10. We also assume, for simplicity, that the cash flow from mining operations is obtained at the *beginning* of the year.

We begin by determining the value of a lease on the mine at time 9, when the remaining deposit is  $x_9$ . Only 1 year remains on the lease, so the value is obtained by maximizing the profit for that year. If we extract  $z_9$  ounces, the revenue from the sale of the gold will be  $gz_9$ , where  $g$  is the price of gold, and the cost of mining will be  $500z_9^2/x_9$ . Hence the optimal value of the mine at time 9 if  $x_9$  is the remaining deposit level is

$$V_9(x_9) = \max_{z_9} (gz_9 - 500z_9^2/x_9).$$

We find the maximum by setting the derivative with respect to  $z_9$  equal to zero. This yields<sup>1</sup>

$$z_9 = gx_9/1,000.$$

We substitute this value in the formula for profit to find

$$V_9(x_9) = \frac{g^2x_9}{1,000} - \frac{500g^2x_9}{1,000 \times 1,000} = \frac{g^2x_9}{2,000}.$$

We write this as  $V_9(x_9) = K_9x_9$ , where  $K_9 = g^2/2,000$  is a constant. Hence the value of the lease is directly proportional to how much gold remains in the mine; the proportionality factor is  $K_9$ .

Next we back up and solve for  $V_8(x_8)$ . In this case we account for the profit generated during the ninth year and also for the value that the lease will have at the end of that year—a value that depends on how much gold we leave in the mine. Hence,

$$V_8(x_8) = \max_{z_8} [gz_8 - 500z_8^2/x_8 + d \cdot V_9(x_8 - z_8)].$$

Note that we have discounted the value associated with the mine at the next year by a factor  $d$ . As in the previous example, the discount rate is constant because the spot rate curve is flat. In this case  $d = 1/1.1$ .

Using the explicit form for the function  $V_9$ , we may write

$$V_8(x_8) = \max_{z_8} [gz_8 - 500z_8^2/x_8 + dK_9(x_8 - z_8)].$$

We again set the derivative with respect to  $z_8$  equal to zero and obtain

$$z_8 = \frac{(g - dK_9)x_8}{1,000}.$$

---

<sup>1</sup> We should check that  $z_9 \leq x_9$ , which does hold with the values we use.

**TABLE 5.4**  
**K-VALUES FOR**  
**COMPLEXICO MINE**

Years	K-values
0	213.81
1	211.45
2	208.17
3	203.58
4	197.13
5	187.96
6	174.79
7	155.47
8	126.28
9	80.00

This value can be substituted into the expression for  $V_8$  to obtain

$$V_8(x_8) = \left[ \frac{(g - dK_9)^2}{2,000} + dK_9 \right] x_8.$$

This is proportional to  $x_8$ , and we may write it as  $V_8(x_8) = K_8 x_8$ .

We can continue backward in this way, determining the functions  $V_7, V_6, \dots, V_0$ . Each of these functions will be of the form  $V_j(x_j) = K_j x_j$ . It should be clear that the same algebra applies at each step, and hence we have the recursive formula

$$K_j = \frac{(g - dK_{j+1})^2}{2,000} + dK_{j+1}.$$

If we use the specific values  $g = 400$  and  $d = 1/1.1$ , we begin the recursion with  $K_9 = g^2/2,000 = 80$ . We can then easily solve for all the other values, as shown in Table 5.4, working from the bottom to the top.

It is the last value calculated (that is,  $K_0$ ) that determines the value of the original lease. That value is determined by finding the value of the lease when there is 50,000 ounces of gold remaining. Hence  $V_0(50,000) = 213.82 \times 50,000 = \$10,691,000$ .

The optimal plan is determined as a by-product of the dynamic programming procedure. At any time  $j$ , the amount of gold to extract is the value  $z_j$  found in the optimization problem. Hence  $z_9 = gx_9/1,000$  and  $z_8 = (g - dK_9)x_8/1,000$ . In general,  $z_j = (g - dK_{j+1})x_j/1,000$ .

Dynamic programming problems using a continuous lattice do not always work out as well as in the preceding example, because it is not always possible to find a simple expression for the  $V$  functions. (The specific functional form for the cost in the gold mine example led to the linear form for the  $V$  functions.) But dynamic programming is a general problem-solving technique that has many variations and many applications. The general idea is used repeatedly in Parts 3 and 4 of this book.

## 5.5 The Harmony Theorem\*

We know that there is a difference between the present value criterion for selecting investment opportunities and the internal rate of return criterion, and that it is strongly believed by theorists that the present value criterion is the better of the two, provided that account is made for the entire cash flow stream of the investment over all its periods. But if you are asked to consider an investment of a fixed amount of dollars (say, in your friend's new venture), you probably would not evaluate this proposition in terms of present value; you would more likely focus on potential return. In fact, if you do make the investment, you are likely to encourage your friend to maximize the return on your investment, not the present value of the firm. Your friend might insist on maximizing present value. Is there a conflict here?

We will try to shed some light on this important issue by working through a hypothetical situation. Suppose your friend has invented a new gismo for which he holds the patent rights. To profit from this invention, he must raise capital and carry out certain operations. The cost for the operations occurs immediately; the reward occurs at the end of a year. In other words, the cash flow stream has just two elements: a negative amount now and a positive amount at the end of a year.

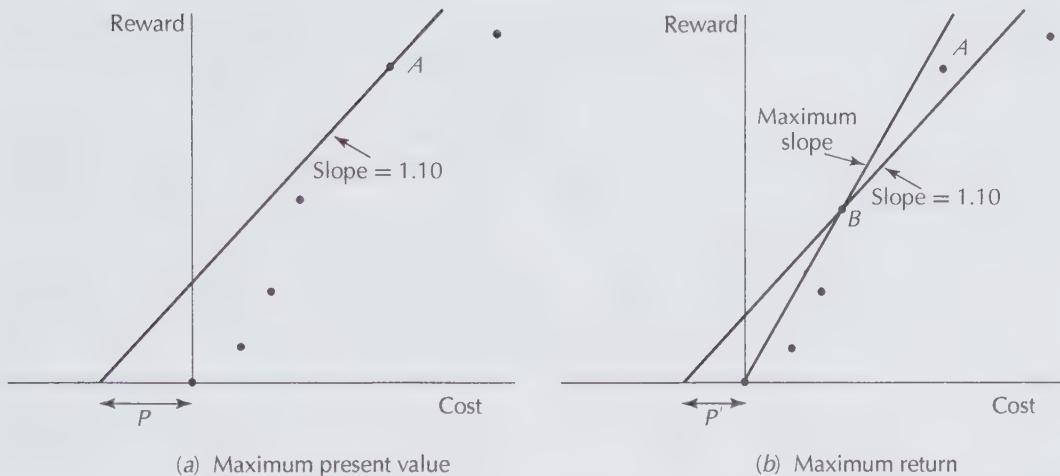
Your friend recognizes that there are many different ways that he can operate his venture, and these entail different costs and different rewards. Hence there are many possible cash flow streams corresponding to different operating plans. He must select one. The possibilities can be described by points on a graph showing the reward (at the end of a year) versus the current cost of operations, as in Figure 5.11(a). Your friend can select any one of the points.

Suppose also that the 1-year interest rate is  $r = 10\%$ . The possibility of depositing money in the bank can be represented on the graph as a straight line with slope 1.10: the current deposit is a cost, and the reward is 1.10 times that amount. This slope will be used to evaluate the present value of a cash flow stream.

If your friend decides to maximize the present value of his venture, he will draw lines with slopes  $1 + r = 1.10$  and find the highest one that goes through a possible operating plan. The plan that lies on that line is the optimal one. This optimal line and plan are shown in Figure 5.11(a); point A is the optimal plan. Using a bank, it is possible to move along the line through A. In particular, it is possible to move all the way down to the horizontal axis. At this point, no money will be received next year, but an amount  $P$  of net profit is obtained now.

Suppose your friend asks you to invest in his venture, supplying a portion of the operating cost and getting that portion of the reward. You would measure the return on your investment. The operating point that achieves the maximum return is found by swinging a line upward, pivoting around the origin, reaching an operating point of greatest possible slope. The result of this process is shown in Figure 5.11(b). The optimal point according to this criterion is the point B in the figure. The maximum return is the slope of this maximum-slope line. Note that this slope is greater than 110%. So point B achieves a higher rate of return than point A. Its present value, however, is just  $P'$ , which is less than  $P$ . There seems to be a conflict.

Here is how the conflict is resolved. Your friend currently owns the rights to his gismo. He has not yet committed any money for operations; but his present value



**FIGURE 5.11 Comparison of criteria.** (a) Plan A is selected because it has the greatest present value. It is the point corresponding to the highest line of slope, equal to 1.10. (b) Plan B is selected because it is the point on the line from the origin of greatest slope. As the text demonstrates, the analysis in (b) is faulty, and when corrected, the maximum return decision will correspond to the present value decision.

analysis shows that it is worth  $P$ . He could go to the bank, take out a loan sufficient to cover the expenses for plan A, and then, at the end of the year, he could pay back the loan and pocket the profit of  $1.10P$  (which is worth  $P$  now). He doesn't care about the rate of return, since he is not investing any money; he is just taking out a loan. Alternatively, he could borrow the money from you, but he would not pay you any more than the current interest rate.

But you are not being asked to make a loan; you are being asked to invest in the venture—to have ownership in it. As an extreme case, suppose your friend asks you to buy the whole venture. You will then have the rights to the gismo. He is willing to stay on and operate the venture (if you provide the necessary operating costs), but you will have the power to decide what operating plan to use.

If your friend sells you the venture, he will charge you an amount  $P$  because that is what it would be worth to him if he kept ownership. So if you decide to buy the venture, the total expense of an operating plan is now  $P$  plus the actual operating cost. If you want to maximize your return, you will maximize  $\text{reward}/(\text{cost} + P)$ . You can find this new best operating plan by swinging a line upward, pivoting around the point  $-P$ , reaching the operating point with the greatest possible slope. That point will be point A, the point that maximized the present value. [Look again at Figure 5.11(a).] Alternatively, once you are the owner, you might consider maximizing the present value. That will lead to point A as well. Therefore if you decide to buy the venture, and you pay the full value  $P$ , you will maximize the return on your investment by operating under plan A; and your return will be 110%. (It does not matter if you decide to borrow some of the operating costs instead of funding them yourself; still you will want to operate at A, and your return will still be 110%.)

We summarize the preceding discussion by a general result that we term the *harmony theorem*. It states that there is harmony between the present value criterion and the rate of return criterion when account is made for ownership.

**Harmony theorem** *Current owners of a venture should want to operate the venture to maximize the present value of its cash flow stream. Potential investors, who must pay the full value of their prospective share of the venture, will want the company to operate in the same way, in order to maximize the return on their investment.*

The harmony theorem is justification for operating a venture (such as a company) in the way that maximizes the present value of the cash flow stream it generates. Both current owners and potential investors will agree on this policy.

The presentation in this section considered only deterministic cash flow streams with two flows. The harmony theorem generalizes to multiple periods and to random streams as well—under certain conditions. A multiperiod generalization is discussed in Exercise 10.

## 5.6 Valuation of a Firm\*

The principles of cash flow analysis can be used to evaluate the worth of publicly traded corporations; indeed almost all analytic valuation methods do use some form of cash flow analysis. However, as straightforward as that may sound, the general idea is subject to a variety of interpretations, each leading to a different result. These differences spring from the question of just which cash flows should form the basis of analysis: should they be the dividends that flow to a stockholder, the net earnings to the company, or the flow that could be captured by a single individual or group who owned the company and was free to extract the cash according to the group's own policy? If these various quantities are defined by standard accounting practice, they can lead to significantly different inferred firm values.

Another weakness of this kind of analysis is that it is based on an assumption that future cash flows are known deterministically, which, of course, is usually not the case. Often uncertainty is recognized in an analysis, but treated in a simplistic way (for instance, by increasing the interest rate used for discounting above the risk-free rate). We discuss other, more solidly based approaches to evaluation under uncertainty in later chapters. This section assumes that the cash flows are deterministic.

### Dividend Discount Models

The owner of a share of stock in a company can expect to receive periodic dividends. Suppose that it is known that in year  $k$ ,  $k = 1, 2, \dots$ , a dividend of  $D_k$  will be received. If the interest rate (or the discount rate) is fixed at  $r$ , it is reasonable to assign a value of the firm to the stock holders as the present value of this dividend stream; namely,

$$V_0 = \frac{D_1}{1+r} + \frac{D_2}{(1+r)^2} + \frac{D_3}{(1+r)^3} + \dots$$

This formula is straightforward, but it requires that the future dividends be known.

A popular way to specify dividends is to use the **constant-growth dividend model**, where dividends grow at a constant rate  $g$ . In particular, given  $D_1$  and the relation  $D_{k+1} = (1 + g)D_k$ , the present value of the stream is

$$V_0 = \frac{D_1}{1+r} + \frac{D_1(1+g)}{(1+r)^2} + \frac{D_1(1+g)^2}{(1+r)^3} + \cdots = D_1 \sum_{k=1}^{\infty} \frac{(1+g)^{k-1}}{(1+r)^k}.$$

This summation is similar to that of an annuity, except that there is the extra growth term in the numerator. The summation will have finite value only if the dividend growth rate is less than the rate used for discounting; that is, if  $g < r$ . In that case we have the explicit **Gordon formula** (see Exercise 11) for the summation

$$V_0 = \frac{D_1}{r - g}. \quad (5.1)$$

Note that, according to this formula, the value of a firm's stock increases if  $g$  increases, if the current dividend  $D_1$  increases, or if the discount rate  $r$  decreases. All of these properties are intuitively clear.

If we project  $D_1$  from a current dividend (already paid) of  $D_0$ , we can rewrite (5.1) by including the first-year's growth. We highlight this as follows:

**Discounted growth formula** Consider a dividend stream that grows at a rate of  $g$  per period. Assign  $r > g$  as the discount rate per period. Then the present value of the stream, starting one period from the present, with the dividend  $D_1$ , is

$$V_0 = \frac{(1+g)D_0}{r - g}, \quad (5.2)$$

where  $D_0$  is the current dividend.

To use the constant-growth dividend model one must estimate the growth rate  $g$  and assign an appropriate value to the discount rate  $r$ . Estimation of  $g$  can be based on the history of the firm's dividends and on future prospects. Frequently a value is assigned to  $r$  that is larger than the actual risk-free interest rate to reflect the idea that uncertain cash flows should be discounted more heavily than certain cash flows. (In Chapters 18 and 19, we study better ways to account for uncertainty.)

**Example 5.6 (The XX Corporation)** The XX Corporation has just paid a dividend of \$1.37M. The company is expected to grow at 10% for the foreseeable future, and hence most analysts project a similar growth in dividends. The discount rate used for this type of company is 15%. What is the value of a share of stock in the XX Corporation?

The total value of all shares is given by (5.2). Hence this value is

$$V_0 = \frac{1.37M \times 1.10}{15 - 10} = \$30,140,000.$$

Assume that there are 1 million shares outstanding. Each share is worth \$30.14 according to this analysis.

## Free Cash Flow\*

A conceptual difficulty with the dividend discount method is that the dividend rate is set by the board of directors of the firm, and this rate may not be representative of the firm's financial status. A different perspective to valuation is obtained by imagining that you were the sole owner and could take out cash. From this perspective the value of the firm might be the discounted value of the net earnings stream.

The net earnings of a firm is defined by accounting practice. In the simplest case it is just revenue minus cost, and then minus taxes; but things are rarely this simple. Account must be made for depreciation of plant and equipment, payment of interest on debt, taxes, and other factors. The final net earnings figure may have little relation to the cash flow that can be extracted from the firm.

Within the limitations of a deterministic approach, the best way to value a firm is to determine the cash flow stream of maximum present value that can be taken out of the company and distributed to the owners. The corresponding cash flow in any year is termed that year's **free cash flow** (FCF). Roughly, free cash flow is the cash generated through operations minus the investments necessary to sustain those operations and their anticipated growth.

It is difficult to obtain an accurate measure of the free cash flow. First, it is necessary to assess the firm's potential for generating cash under various policies. Second, it is necessary to determine the optimal rate of investment—the rate that will generate the cash flow stream of maximum present value. Usually this optimal rate is merely estimated; but since the relation between growth rate and present value is complex, the estimated rate may be far from the true optimum. We shall illustrate the ideal process with a highly idealized example.

Suppose that a company has gross earnings of  $Y_n$  in year  $n$  and decides to invest a portion  $u$  of this amount each year in order to attain earnings growth. The growth rate is determined by the function  $g(u)$ , which is a property of the firm's characteristics. On a (simplified) accounting basis, depreciation is a fraction  $\alpha$  of the current capital account ( $\alpha \approx .10$ , for example). In this case the capital  $C_n$  follows the formula  $C_{n+1} = (1 - \alpha)C_n + uY_n$ . With these ideas we can set up a general income statement for a firm, as shown in Table 5.5.

**Example 5.7 (Optimal growth)** We can go further with the foregoing analysis and calculate  $Y_n$  and  $C_n$  in explicit form. Since  $Y_{n+1} = [1 + g(u)]Y_n$ , it is easy to see that  $Y_n = [1 + g(u)]^n Y_0$ . Likewise, it can be shown that

$$C_n = (1 - \alpha)^n C_0 + uY_0 \left\{ \frac{-(1 - \alpha)^n + [1 + g(u)]^n}{g(u) + \alpha} \right\}.$$

If we ignore the two terms having  $(1 - \alpha)^n$  (since they go to zero for large  $n$ ) we have

$$C_n = \frac{uY_0[1 + g(u)]^n}{g(u) + \alpha}. \quad (5.3)$$

**TABLE 5.5**  
**FREE CASH FLOW**

Income statement	
Before-tax cash flow from operations	$Y_n$
Depreciation	$\alpha C_n$
Taxable income	$Y_n - \alpha C_n$
Taxes (34%)	$.34(Y_n - \alpha C_n)$
After-tax income	$.66(Y_n - \alpha C_n)$
After-tax cash flow (after-tax income plus depreciation)	$.66(Y_n - \alpha C_n) + \alpha C_n$
Sustaining investment	$u Y_n$
Free cash flow	$.66(Y_n - \alpha C_n) + \alpha C_n - u Y_n$

Depreciation is assumed to be  $\alpha$  times the amount in the capital account.

Putting the expressions for  $Y_n$  and  $C_n$  in the bottom line of Table 5.5, we find the free cash flow at time  $n$  to be

$$FCF = \left[ .66 + .34 \frac{\alpha u}{g(u) + \alpha} - u \right] [1 + g(u)]^n Y_0. \quad (5.4)$$

This is a growing geometric series. We can use the Gordon formula to calculate its present value at interest rate  $r$ . This gives

$$PV = \left[ .66 + .34 \frac{\alpha u}{g(u) + \alpha} - u \right] \frac{1+r}{r-g(u)} Y_0. \quad (5.5)$$

It is not easy to see by inspection what value of  $u$  would be best. Let us consider another example.

**Example 5.8 (XX Corporation)** Assume that the XX Corporation has current earnings of  $Y_0 = \$10$  million, and the initial capital<sup>2</sup> is  $C_0 = \$19.8$  million. The interest rate is  $r = 15\%$ , the depreciation factor is  $\alpha = .10$ , and the relation between investment rate and growth rate is  $g(u) = .12[1 - e^{5(\alpha-u)}]$ . Notice that  $g(\alpha) = 0$ , reflecting the fact that an investment rate of  $\alpha$  times earnings just keeps up with the depreciation of capital.

Using (5.5) we can find the value of the company for various choices of the investment rate  $u$ . For example, for  $u = 0$ , no investment, the company will slowly shrink, and the present value under that policy will be \$33.4 million. If  $u = .10$ , the company will just maintain its current level, and the present value under that plan will be \$45.5 million. Or if  $u = .5$ , the present value will be \$59.8 million.

It is possible to maximize (5.5) (by trial and error or by a simple optimization routine as is available in some spreadsheet packages). The result is  $u = 37.7\%$  and  $g(u) = 9.0\%$ . The corresponding present value is \$67.0 million. This is the company value.

<sup>2</sup> This value of  $C_0$  will make the terms that were canceled in deriving (5.3) cancel exactly.

Here is a question to consider carefully. Suppose that during the first year, the firm operates according to this plan, investing 37.7% of its gross earnings in new capital. Suppose also, for simplicity, that no dividends are paid that year. What will be the value of the company after 1 year? Recall that during this year, capital and earnings expand by 9%. Would you guess that the company value will increase by 9% as well? Remember the harmony theorem. Actually, the value will increase by the rate of interest, which is 15%. Investors must receive this rate, and they do. The reason this may seem strange is that we assumed that no dividends were paid. The free cash flow that was generated, but not taken out of the company, is held for the year (itself earning 15%), and this must be added to the present value calculation of future cash flows. If the free cash flow generated in the first year were distributed as dividends, the company value would increase by 9%, but the total return to investors, including the dividend and the value increase, again would be 15%.

Although this example is highly idealized, it indicates the character of a full valuation procedure (under an assumption of certainty). The free cash flow stream must be projected, accounting for future opportunities. Furthermore, this cash flow stream must be optimized by proper selection of a capital investment policy. Because the impact of current investment on future free cash flow is complex, effective optimization requires the use of formal models and formal optimization techniques.

## 5.7 Summary

Interest rate theory is probably the most widely used financial tool. It is used to determine the value of projects, to allocate money among alternatives, to design complex bond portfolios, to determine how to manage investments effectively, and even to determine the value of a firm.

Interest rate theory is most powerful when it is combined with general problem-solving methods, particularly methods of optimization. With the aid of such methods, interest rate theory provides more than just a static measure of value; it guides us to find the decision or structure with the highest value.

One class of problems that can be approached with this combination is capital budgeting problems. In the classic problem of this class, a fixed budget is to be allocated among a set of independent projects in order to maximize net present value. This problem can be solved approximately by selecting projects with the highest benefit-cost ratio. The problem can be solved exactly by formulating it as a zero-one optimization problem and using an integer programming package. More complex capital budgeting problems having dependencies among projects can be also be solved by the zero-one programming method.

The selection of a bond portfolio to meet certain requirements can be conveniently formulated as an optimization problem—but there are several possible formulations. A particularly simple problem within this class is the cash-matching problem, where a portfolio is constructed to generate a required cash flow in each

period. This formulation has the weakness that in some periods extra cash may be generated, beyond that required, and this extra cash is essentially wasted. More complex formulations do not have this weakness.

To produce excellent results, many investments require deliberate ongoing management. The relation between a series of management decisions and the resulting cash flow stream frequently can be modeled as a graph. (Especially useful types of graphs are trees and lattices.) In such a graph the nodes correspond to states of the process, and a branch leading from a node corresponds to a particular choice made from that node. Associated with each branch is a cash flow value.

Optimal dynamic management consists of following the special path of arcs through the graph that produces the greatest present value. This optimal path can be found efficiently by the method of dynamic programming. A particularly useful version of dynamic programming for investment problems uses the running method for evaluation of present value.

Dynamic programming works backward in time. For a problem with  $n$  time periods, the running version of the procedure starts by finding the best decision at each of the nodes  $i$  at time  $n - 1$  and assigns a  $V$ -value, denoted by  $V_{n-1,i}$ , to each such node. This  $V$ -value is the optimal present value that could be obtained if the investment process were initiated at that node. To find that value, each possible arc emanating from node  $i$  is examined. The sum of the cash flow of the arc and the one-period discounted  $V$ -value at the node reached by the arc is evaluated. The  $V$ -value of the originating node  $i$  is the maximum of those sums. After completing this procedure for all the nodes at  $n - 1$ , the procedure then steps back to the nodes at time  $n - 2$ . Optimal  $V$ -values are found for each of those nodes by a procedure that exactly parallels that for the nodes at  $n - 1$ . The procedure continues by working backward through all time periods, and it ends when an optimal  $V$ -value is assigned to the initial node at time zero.

When operating a venture it is appropriate to maximize the present value. On the other hand, investors may be most interested in the rate of return. These criteria might seem to be in conflict, but the **harmony theorem** states that the criteria are equivalent under the assumption that investors pay the full value for their ownership of the venture.

Present value analysis is commonly used to estimate the value of a firm. One such procedure is the dividend discount method, where the value to a stockholder is assumed to be equal to the present value of the stream of future dividend payments. If dividends are assumed to grow at a rate  $g$  per year, a simple formula gives the present value of the resulting stream.

The better method of firm evaluation bases the evaluation on free cash flow, which is the amount of cash that can be taken out of the firm while maintaining optimal operations and investment strategies. In idealized form, this method requires that the present value of free cash flow be maximized with respect to all possible management decisions, especially those related to investment that produces earnings growth.

Valuation methods based on present value suffer the defect that future cash flows are treated as if they were known with certainty, when in fact they are usually uncertain. The deterministic theory is therefore not adequate. This defect is widely recognized; and to compensate for it, it is common practice to discount predicted,

but uncertain, cash flows at higher interest rates than the risk-free rate. There is some theoretical justification for this, but a completely consistent approach to uncertainty is more subtle. The exciting story of uncertainty in investment begins with the next chapter and continues throughout the remainder of the text.

## Exercises

- (Capital budgeting) A firm is considering funding several proposed projects that have the financial properties shown in Table 5.6. The available budget is \$600,000. What set of projects would be recommended by the approximate method based on benefit–cost ratios? What is the optimal set of projects?

**TABLE 5.6**  
**FINANCIAL PROPERTIES OF PROPOSED PROJECTS**

Project	Outlay (\$1,000)	Present worth (\$1,000)
1	100	200
2	300	500
3	200	300
4	150	200
5	150	250

- (The road ⊕) Refer to the transportation alternatives problem of Example 5.2. The bridge at Cay Road is actually part of the road between Augen and Burger. Therefore it is not reasonable for the bridge to have fewer lanes than the road itself. This means that if projects 2 or 4 are carried out, either projects 6 or 7 must also be carried out. Formulate a zero–one programming problem that includes this additional requirement. Solve the problem.
- (Two-period budget ⊕) A company has identified a number of promising projects, as indicated in Table 5.7. The cash flows for the first 2 years are shown (they are all negative). The cash flows in later years are positive, and the net present value of each project is shown. The company managers have decided that they can allocate up to \$250,000 in each of the first 2 years to fund these projects. If less than \$250,000 is used the first year, the balance

**TABLE 5.7**  
**A LIST OF PROJECTS**

Project	Cash flow		NPV
	1	2	
1	−90	−58	150
2	−80	−80	200
3	−50	−100	100
4	−20	−64	100
5	−40	−50	120
6	−80	−20	150
7	−80	−100	240

can be invested at 10% and used to augment the next year's budget. Which projects should be funded?

4. (Bond matrix  $\diamond$ ) The cash matching and other problems can be conveniently represented in matrix form. Suppose there are  $m$  bonds. We define for each bond  $j$  its associated yearly cash flow stream (column) vector  $\mathbf{c}_j$ , which is  $n$ -dimensional. The yearly obligations are likewise represented by the  $n$ -dimensional vector  $\mathbf{y}$ . We can stack the  $\mathbf{c}_j$  vectors side by side to form the columns of a bond matrix  $\mathbf{C}$ . Finally we let  $\mathbf{p}$  and  $\mathbf{x}$  be  $m$ -dimensional column vectors. The cash matching problem can be expressed as

$$\begin{aligned} & \text{maximize} && \mathbf{p}^T \mathbf{x} \\ & \text{subject to} && \mathbf{C}\mathbf{x} \geq \mathbf{y} \\ & && \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

(a) Identify  $\mathbf{C}, \mathbf{y}, \mathbf{p}$ , and  $\mathbf{x}$  in Table 5.3.

(b) Show that if all bonds are priced according to a common term structure of interest rates, there is a vector  $\mathbf{v}$  satisfying

$$\mathbf{C}^T \mathbf{v} = \mathbf{p}.$$

What are the components of  $\mathbf{v}$ ?

(c) Suppose  $\mathbf{b}$  is a vector whose components represent obligations in each period. Show that a portfolio  $\mathbf{x}$  meeting these obligations exactly satisfies

$$\mathbf{C}\mathbf{x} = \mathbf{b}.$$

(d) With  $\mathbf{x}$  and  $\mathbf{v}$  defined as before, show that the price of the portfolio  $\mathbf{x}$  is  $\mathbf{v}^T \mathbf{b}$ . Interpret this result.

5. (Trinomial lattice) A trinomial lattice is a special case of a trinomial tree. From each node three moves are possible: up, middle, and down. The special feature of the lattice is that certain pairs of moves lead to identical nodes two periods in the future. We can express these equivalences as

$$\text{up-down} = \text{down-up} = \text{middle-middle}$$

$$\text{middle-down} = \text{down-middle}$$

$$\text{middle-up} = \text{up-middle}.$$

Draw a trinomial lattice spanning three periods. How many nodes does it contain? How many nodes are contained in a full trinomial tree of the same number of periods?

6. (A bond project  $\oplus$ ) You are the manager of XYZ Pension Fund. On November 5, 2021, XYZ must purchase a portfolio of U.S. Treasury bonds to meet the fund's projected liabilities in the future. The bonds available at that time are those of Exercise 4 in Chapter 4. Short selling is not allowed. Following the procedure of the earlier exercise, a 4th-order polynomial estimate of the term structure is constructed as  $r(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \alpha_3 t^3 + \alpha_4 t^4$ . The liabilities of XYZ are as listed in Table 5.8.

(a) (Simple cash matching) Construct a minimum-cost liability-matching portfolio by buying Treasury bonds assuming that excess periodic cash flows may be held only at zero interest to meet future liabilities.

**TABLE 5.8**  
**LIABILITIES OF XYZ**  
**PENSION FUND**

Liabilities	Occur on 15th
Feb 2022	\$2,000
Aug 2022	\$20,000
Feb 2023	\$0
Aug 2023	\$25,000
Feb 2024	\$1,000
Aug 2024	\$0
Feb 2025	\$20,000
Aug 2025	\$1,000
Feb 2026	\$15,000

- (b) (Complex cash matching) Construct a minimum-cost liability-matching portfolio by buying Treasury bonds assuming that all excess periodic cash flows may be reinvested at the expected interest rates (implied by the current term structure) to meet future liabilities. No borrowing is allowed.
- (c) (Duration matching) Construct a minimum-cost portfolio with present value equal to that of the liability stream. Immunize against a change in the term structure parameters. Do this for five cases. Case 1 is to guard against a change in  $\alpha_1$ , case 2 to guard against changes in  $\alpha_1$  and  $\alpha_2$ , and so on.
7. (The fishing problem) Find the solution to the fishing problem of Example 5.4 when the interest rate is 33%. Are the decisions different than when the interest rate is 25%? At what critical value of the discount factor does the solution change?
8. (Complexico mine  $\oplus$ ) Consider the Complexico mine and assume a 10% constant interest rate; also assume the price of gold is constant at \$400/oz.
- (a) Find the value of the mine (not a 10-year lease) if the current deposit is  $x_0$ . In particular, how much is the mine worth initially when  $x_0 = 50,000$  ounces? [Hint: Consider the recursive equation for  $K_k$  as  $k \rightarrow \infty$ .]
- (b) For the 10-year lease considered in the text, how much gold remains in the mine at the end of the lease; and how much is the mine worth at that time?
- (c) If the mine were not leased, but instead operated optimally by an owner, what would the mine be worth after 10 years?
9. (Little Bear Oil) You have purchased a lease for the Little Bear Oil well. This well has initial reserves of 100 thousand barrels of oil. In any year you have three choices of how to operate the well: (a) you can *not* pump, in which case there is no operating cost and no change in oil reserves; (b) you can pump normally, in which case the operating cost is \$50 thousand and you will pump out 20% of what the reserves were at the beginning of the year; or (c) you can use enhanced pumping using water pressure, in which case the operating cost is \$120 thousand and you will pump out 36% of what the reserves were at the beginning of the year. The price of oil is \$10 per barrel and the interest rate is 10%. Assume that both your operating costs and the oil revenues come at the beginning of the year (through advance sales). Your lease is for a period of 3 years.

- (a) Show how to set up a trinomial lattice to represent the possible states of the oil reserves.
- (b) What is the maximum present value of your profits, and what is the corresponding optimal pumping strategy?
- 10.** (Multiperiod harmony theorem  $\diamond$ ) The value of a firm is the maximum present value of its possible cash flow streams. This can be expressed as

$$V_0 = \max \left[ x_0 + \frac{x_1}{1+s_1} + \frac{x_2}{(1+s_2)^2} + \cdots + \frac{x_n}{(1+s_n)^n} \right],$$

where the maximization is with respect to all possible streams  $x_0, x_1, \dots, x_n$ , and the  $s_i$ 's are the spot rates. Let  $x_0^*$  be the first cash flow in the optimal plan. If the firm chooses an arbitrary plan that results in an initial cash flow of  $x_0$  (distributed to the owners), the value of the firm after 1 year is

$$V_1(x_0) = \max \left\{ x_1 + \frac{x_2}{1+s'_1} + \frac{x_3}{(1+s'_2)^2} + \cdots + \frac{x_n}{(1+s'_n)^{n-1}} \right\},$$

where now that maximum is with respect to all feasible cash flows that start with  $x_0$  and the  $s'_i$ 's are the spot rates after 1 year. An investor purchasing the firm at its full fair price has initial cash flow  $x_0 - V_0$  and achieves a value of  $V_1(x_0)$  after 1 year. Hence the 1-year total return to the investor is

$$R = \frac{V_1(x_0)}{V_0 - x_0}.$$

The investor would urge that  $x_0$  be chosen to maximize  $R$ . Call this value  $\bar{x}_0$ . Assuming that interest rates follow expectation dynamics and that  $V_1(\bar{x}_0) > 0$ , show that the maximum  $R$  is  $1 + s_1$  and that this return is achieved by the same  $x_0^*$  that determines  $V_0$ .

- 11.** (Growing annuity) Show that for  $g < r$ ,

$$\sum_{k=1}^{\infty} \frac{(1+g)^{k-1}}{(1+r)^k} = \frac{1}{r-g}.$$

[Hint: Let  $S$  be the value of the sum. Note that  $S = 1/(1+r) + S(1+g)/(1+r)$ .]

- 12.** (Two-stage growth) It is common practice in security analysis to modify the basic dividend growth model by allowing more than one stage of growth, with the growth factors being different in the different stages. As an example consider company Z, which currently distributes dividends of \$10M annually. The dividends are expected to grow at the rate of 10% for the next 5 years and at a rate of 5% thereafter.

- (a) Using a dividend discount approach with an interest rate of 15%, what is the value of the company?
- (b) Find a general formula for the value of a company satisfying a two-stage growth model. Assume a growth rate of  $G$  for  $k$  years, followed by a growth rate of  $g$  thereafter, and an initial dividend of  $D_1$ .

## References

Capital budgeting is a classic topic in financial planning. Some good texts are [1–4]. Bond portfolio construction is considered in [5, 6] and in other references given for Chapters 3 and 4. Dynamic programming was developed by Bellman (see [7, 8]). The Harmony Theorem was stated in [9]. The

classic reference on stock valuation is [10]. See [11–14] for other presentations. A vivid discussion of how improper analysis techniques led to disastrous overvaluation in the 1980s is in [15].

1. Peterson, P. P., and Fabozzi, F. J. (2002), *Capital Budgeting*, John Wiley, & Sons, New York.
2. Brealey, R., and S. Myers (2010), *Principles of Corporate Finance*, 10th ed., McGraw-Hill, New York.
3. Bierman, H., Jr., and S. Smidt (2006), *The Capital Budgeting Decision*, 9th ed., Macmillan, New York.
4. Baker, H. K. (2011), *Capital Budgeting Valuation*, John Wiley & Sons, New York.
5. Bierwag, G. O., G. G. Kaufman, R. Schweitzer, and A. Toebs (1981), “The Art of Risk Management in Bond Portfolios,” *Journal of Portfolio Management*, **7**, 27–36.
6. Fabozzi, F. J., and T. D. Fabozzi (1989), *Bond Markets, Analysis and Strategies*, Prentice Hall, Englewood Cliffs, NJ.
7. Bellman, R. (1957), *Dynamic Programming*, Princeton University Press, Princeton, NJ.
8. Bellman R., and S. Dreyfus (1962), *Applied Dynamic Programming*, Princeton University Press, Princeton, NJ.
9. Luenberger, D. G. (1998), *Investment Science*, 1st ed., Oxford University Press, New York.
10. Graham, B., D. L. Dodd, and S. Cottle (1962), *Security Analysis*, McGraw-Hill, New York.
11. Williams, J. B. (1938), *The Theory of Investment Value*, North-Holland, Amsterdam, The Netherlands.
12. Gordon, M. J. (1959), “Dividends, Earnings, and Stock Prices,” *Review of Economics and Statistics*, **41**, 99–195.
13. Fridson, M. S., and F. Alvarez (2011), *Financial Statement Analysis*, John Wiley & Sons, New York.
14. Black, F. (1980), “The Magic in Earnings: Economic Earnings versus Accounting Earnings,” *Financial Analysts Journal*, **36**, 19–24.
15. Klarman, S. A. (1991), *Margin of Safety: Risk-Averse Value Investing Strategies for the Thoughtful Investor*, Harper Business.

## PART II

# SINGLE-PERIOD RANDOM CASH FLOWS





# MEAN-VARIANCE PORTFOLIO THEORY

**T**ypically, when making an investment, the initial outlay of capital is known, but the amount to be returned is uncertain. Such situations are studied in this part of the text. In this part, however, we restrict attention to the case of a single investment period: money is invested at the initial time, and payoff is attained at the end of the period.

The assumption that an investment situation comprises a single period is sometimes a good approximation. An investment in a zero-coupon bond that will be held to maturity is an example. Another is an investment in a physical project that will not provide payment until it is completed. However, many common investments, such as publicly traded stocks, are not tied to a single period, since they can be liquidated at will and may return dividends periodically. Nevertheless, such investments are often analyzed on a single period basis as a simplification; but this type of analysis should be regarded only as a prelude to Parts 3 and 4 of the text, which are more comprehensive.

This part of the text treats uncertainty with three different mathematical methods: (1) mean-variance analysis, (2) utility function analysis, and (3) arbitrage (or comparison) analysis. Each of these methods is an important component of investment science.

This first chapter of the second part of the text treats uncertainty by **mean-variance** analysis. This method uses probability theory only slightly, and leads to convenient mathematical expressions and procedures. Mean-variance analysis forms the basis for the important *capital asset pricing model* discussed in Chapter 7.

## 6.1 Asset Return

An investment instrument that can be bought and sold is frequently called an **asset**. We introduce a fundamental concept concerning such assets.

Suppose that you purchase an asset at time zero, and 1 year later you sell the asset. The **total return** on your investment is defined to be

$$\text{total return} = \frac{\text{amount received}}{\text{amount invested}}.$$

Or if  $X_0$  and  $X_1$  are, respectively, the amounts of money invested and received and  $R$  is the total return, then

$$R = \frac{X_1}{X_0}.$$

Often, for simplicity, the term *return* is used for total return.

The **rate of return** is

$$\text{rate of return} = \frac{\text{amount received} - \text{amount invested}}{\text{amount invested}}.$$

Or, again, if  $X_0$  and  $X_1$  are, respectively, the amounts of money invested and received and  $r$  is the rate of return, then

$$r = \frac{X_1 - X_0}{X_0}. \quad (6.1)$$

The shorter expression *return* is also frequently used for the rate of return.

We distinguish the two definitions by using upper- or lowercase letters, such as  $R$  and  $r$ , respectively, for total return and rate of return; and usually the context makes things clear if we use the shorthand phrase *return*.

It is clear that the two notions are related by

$$R = 1 + r$$

and that equation (6.1) can be rewritten as

$$X_1 = (1 + r)X_0.$$

This shows that a rate of return acts much like an interest rate.

## Short Sales

Sometimes it is possible to sell an asset that you do not own through the process of **short selling**, or **shorting**, the asset. To do this, you borrow the asset from someone who owns it (such as a brokerage firm). You then sell the borrowed asset to someone else, receiving an amount  $X_0$ . At a later date, you repay your loan by purchasing the asset for, say,  $X_1$  and return the asset to your lender. If the later amount  $X_1$  is lower than the original amount  $X_0$ , you will have made a profit of  $X_0 - X_1$ . Hence short selling is profitable if the asset price declines.

Short selling is considered quite risky—even dangerous—by many investors. The reason is that the potential for loss is unlimited. If the asset value increases, the loss is  $X_1 - X_0$ ; since  $X_1$  can increase arbitrarily, so can the loss. For this reason (and others) short selling is prohibited within certain financial institutions, and it is purposely avoided as a policy by many individuals and institutions. However, it is not universally forbidden, and there is, in fact, a considerable level of short selling of stock market securities.

When short selling a stock, you are essentially duplicating the role of the issuing corporation. You sell the stock to raise immediate capital. If the stock pays dividends during the period that you have borrowed it, you too must pay that same dividend to the person from whom you borrowed the stock.

In practice, the pure process of short selling is supplemented by certain restrictions and safeguards. (For example, you must post a security deposit with the broker from whom you borrowed the asset.) But for theoretical work, we typically assume that the pure shorting of an asset is allowed.

Let us determine the return associated with short selling. We receive  $X_0$  initially and pay  $X_1$  later, so the outlay is  $-X_0$  and the final receipt is  $-X_1$ , and hence the total return is

$$R = \frac{-X_1}{-X_0} = \frac{X_1}{X_0}.$$

The minus signs cancel out, so we obtain the same expression as that for purchasing the asset. Hence the return value  $R$  applies algebraically to both purchases and short sales. We can write this as

$$-X_1 = -X_0 R = -X_0(1 + r)$$

to show that final receipt is related to initial outlay.

**Example 6.1 (A short sale)** Suppose I decide to short 100 shares of stock in company CBA. This stock is currently selling for \$10 per share. I borrow 100 shares from my broker and sell these in the stock market, receiving \$1,000. At the end of 1 year the price of CBA has dropped to \$9 per share. I buy back 100 shares for \$900 and give these shares to my broker to repay the original loan. Because the stock price fell, this has been a favorable transaction for me. I made a profit of \$100.

Someone who purchased the stock at the beginning of the year and sold it at the end would have lost \$100. That person would easily compute

$$R = \frac{900}{1,000} = .90$$

or

$$r = \frac{900 - 1,000}{1,000} = -.10.$$

The rate of return is clearly negative as  $r = -10\%$ . Shorting converts a negative rate of return into a profit because the original investment is also negative. For my shorting activity on CBA my original outlay was  $-\$1,000$ ; hence my profit is  $-\$1,000 \times r = \$100$ .

It is a bit strange to refer to a rate of return associated with the idealized shorting procedure, since there is no initial commitment of resources. Nevertheless, it is the proper notion. In practice, shorting does require an initial commitment of margin, and the proceeds from the initial sale are held until the short is cleared. This modified procedure will have a different rate of return. (See Exercise 1.) For basic theoretical work, however, we shall often assume that the idealized procedure is available.

## Portfolio Return

Suppose now that  $n$  different assets are available. We can form a **master asset**, or **portfolio**, of these  $n$  assets. Suppose that this is done by apportioning an amount  $X_0$  among the  $n$  assets. We then select amounts  $X_{0i}$ ,  $i = 1, 2, \dots, n$ , such that  $\sum_{i=1}^n X_{0i} = X_0$ , where  $X_{0i}$  represents the amount invested in the  $i$ th asset. If we are allowed to sell an asset short, then some of the  $X_{0i}$ 's can be negative; otherwise we restrict the  $X_{0i}$ 's to be nonnegative.

The amounts invested can be expressed as fractions of the total investment. Thus we write

$$X_{0i} = w_i X_0, \quad i = 1, 2, \dots, n,$$

where  $w_i$  is the **weight** or fraction of asset  $i$  in the portfolio. Clearly,

$$\sum_{i=1}^n w_i = 1$$

and some  $w_i$ 's may be negative if short selling is allowed.

Let  $R_i$  denote the total return of asset  $i$ . Then the amount of money generated at the end of the period by the  $i$ th asset is  $R_i X_{0i} = R_i w_i X_0$ . The total amount received by this portfolio at the end of the period is therefore  $\sum_{i=1}^n R_i w_i X_0$ . Hence we find that the overall total return of the portfolio is

$$R = \frac{\sum_{i=1}^n R_i w_i X_0}{X_0} = \sum_{i=1}^n w_i R_i.$$

Equivalently, since  $\sum_{i=1}^n w_i = 1$ , we have

$$r = \sum_{i=1}^n w_i r_i.$$

This is a basic result concerning returns, and so we highlight it here:

**Portfolio return** *Both the total return and the rate of return of a portfolio of assets are equal to the weighted sum of the corresponding individual asset returns, with the weight of an asset being its relative weight (in purchase cost) in the portfolio; that is,*

$$R = \sum_{i=1}^n w_i R_i, \quad r = \sum_{i=1}^n w_i r_i.$$

**TABLE 6.1**  
**CALCULATION OF PORTFOLIO RETURN**

Security	Number of shares	Price	Total cost	Weight in portfolio
Jazz, Inc.	100	\$40	\$4,000	0.25
Classical, Inc.	400	\$20	\$8,000	0.50
Rock, Inc.	200	\$20	\$4,000	0.25
Portfolio total values			\$16,000	1.00
Security	Weight in portfolio	Rate of return	Weighted rate	
Jazz, Inc.	.25	17%	4.25%	
Classical, Inc.	.50	13%	6.50%	
Rock, Inc.	.25	23%	5.75%	
Portfolio rate of return			16.50%	

The weight of a security in a portfolio is its proportion of total cost, as shown in the upper table. These weights then determine the rate of return of the portfolio, as shown in the lower table.

An example calculation of portfolio weights and the associated expected rate of return of the portfolio are shown in Table 6.1.

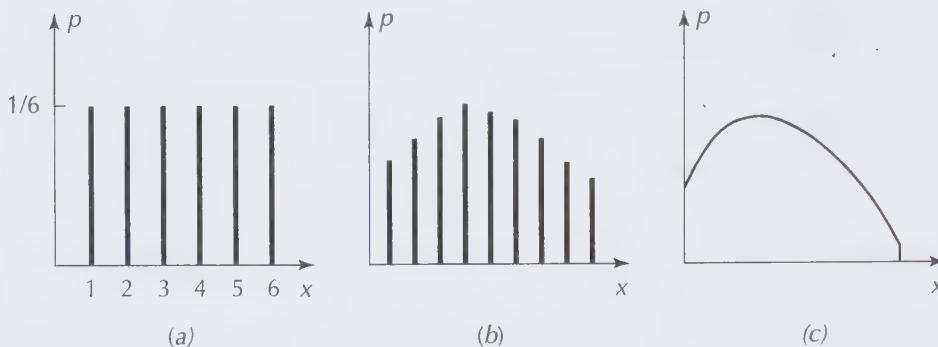
## 6.2 Random Variables

Frequently the amount of money to be obtained when selling an asset is uncertain at the time of purchase. In that case the return is random and can be described in probabilistic terms. In preparation for the study of random returns, we briefly introduce some concepts of probability. (For more detail on basic probability theory, see Appendix A.) Readers with a basic knowledge of probability may wish to turn directly to the next section.

Suppose  $x$  is a random quantity that can take on any one of a finite number of specific values, say,  $x_1, x_2, \dots, x_m$ . Assume further that associated with each possible  $x_i$ , there is a probability  $p_i$  that represents the relative chance of an occurrence of  $x_i$ . The  $p_i$ 's satisfy  $\sum_{i=1}^m p_i = 1$  and  $p_i \geq 0$  for each  $i$ . Each  $p_i$  can be thought of as the relative frequency with which  $x_i$  would occur if an experiment of observing  $x$  were repeated infinitely often. The quantity  $x$ , characterized in this way before its value is known, is called a **random variable**.

A simple example is that of rolling an ordinary six-sided die, with the number of spots obtained being  $x$ . The six possibilities are 1, 2, 3, 4, 5, 6, and each has probability 1/6.

It is common to display the probabilities associated with a random variable graphically as a density. The possible values of  $x$  are indicated on the horizontal axis, and the height of the line at a point represents the probability of that point. Some



**FIGURE 6.1 Probabilities.** Probabilities are shown for (a) the outcome of a roll of a die, (b) another random variable with a finite number of possible outcomes, and (c) a continuous random variable.

examples are shown in Figure 6.1. Figure 6.1(a) shows the probabilities corresponding to the outcome of a roll of a die, where the six possibilities each have a probability of  $1/6$ . Figure 6.1(b) shows a more general case with several possible outcomes of various probabilities.

If the outcome variable can take any real value in an interval as, for example, the temperature of a room, a **probability density function**  $p(x)$  is used. Roughly,  $p(x)dx$  is the probability that the random variable has a value in the interval  $[x, x+dx]$ . The probability that the variable's value will lie in any segment of the line is equal to the area of the vertical region bounded by this segment and the density function. An example is shown in Figure 6.1(c).

## Expected Value

The expected value of a random variable  $x$  is just the average value obtained by regarding the probabilities as frequencies. For the case of a finite number of possibilities, it is defined as

$$E(x) = \sum_{i=1}^m x_i p_i.$$

For convenience  $E(x)$  is often denoted by  $\bar{x}$ . Also the terms **mean** or **mean value** are often used for the expected value. So we say  $x$  has mean  $\bar{x}$ .

**Example 6.2 (A roll of the die)** The expected value of the number of spots on a roll of a die is

$$\frac{1}{6}(1 + 2 + 3 + 4 + 5 + 6) = 3.5.$$

Note that the expected value is not necessarily a possible outcome of a roll.

The expected value operation  $E$  is the main operation used in probability calculations, so it is useful to note its basic properties:

- 1. Certain value** If  $y$  is a known value (not random), then  $E(y) = y$ .

This states that the expected value of a nonrandom quantity is equal to the quantity itself.

- 2. Linearity** If  $y$  and  $z$  are random, then  $E(\alpha y + \beta z) = \alpha E(y) + \beta E(z)$  for any real values of  $\alpha$  and  $\beta$ .

This states that the expected (or mean) value of the sum of two random variables is the sum of their corresponding means; and the mean value of the multiple of a random variable is the same multiple of the original mean. For example, the expected value for the total number of spots on two dice is  $3.5 + 3.5 = 7$ .

- 3. Nonnegativity** If  $x$  is random but never less than zero, then  $E(x) \geq 0$ .

This is a sign-preserving property.

## Variance

The expected value of a random variable provides a useful summary of the probabilistic nature of the variable. However, typically one wants, in addition, to have a measure of the degree of possible deviation from the mean. One such measure is the **variance**.

Given a random variable  $y$  with expected value  $\bar{y}$ , the quantity  $y - \bar{y}$  is itself random, but has an expected value of zero. [This is because  $E(y - \bar{y}) = E(y) - E(\bar{y}) = \bar{y} - \bar{y} = 0$ .] The quantity  $(y - \bar{y})^2$  is always nonnegative and is large when  $y$  deviates greatly from  $\bar{y}$  and small when it is near  $\bar{y}$ . The expected value of this squared variable  $(y - \bar{y})^2$  is a useful measure of how much  $y$  tends to vary from its expected value.

In general, for any random variable  $y$  the variance of  $y$  is defined as

$$\text{var}(y) = E[(y - \bar{y})^2].$$

In mathematical expressions, variance is represented by the symbol  $\sigma^2$ . Thus we write  $\sigma_y^2 = \text{var}(y)$ , or if  $y$  is understood, we simply write  $\sigma^2 = \text{var}(y)$ .

We frequently use the square root of the variance, denoted by  $\sigma$  and called the **standard deviation**. It has the same units as the quantity  $y$  and is another measure of how much the variable is likely to deviate from its expected value. Thus, formally,

$$\sigma_y = \sqrt{E[(y - \bar{y})^2]}.$$

There is a simple formula for variance that is useful in computations. We note that

$$\begin{aligned}\text{var}(x) &= E[x] \\ &= E(x^2) - 2E(x)\bar{x} + \bar{x}^2 \\ &= E(x^2) - \bar{x}^2.\end{aligned}\tag{6.2}$$

This result is used in the following example.

**Example 6.3 (A roll of the die)** Let us compute the variance of the random variable  $y$  defined as the number of spots obtained by a roll of a die. Recalling that  $\bar{y} = 3.5$  we find

$$\begin{aligned}\sigma^2 &= E(y^2) - \bar{y}^2 \\ &= \frac{1}{6}[1 + 4 + 9 + 16 + 25 + 36] - (3.5)^2 = 2.92.\end{aligned}$$

Hence  $\sigma = \sqrt{2.92} = 1.71$ .

## Several Random Variables

Suppose we are interested in two random variables, such as the outside temperature and the barometric pressure. To describe these random variables we must have probabilities for all possible combinations of the two values. If we denote the variables by  $x$  and  $y$ , we must consider the possible pairs  $(x, y)$ . Suppose  $x$  can take on the possible values  $x_1, x_2, \dots, x_n$  and  $y$  can take on the values  $y_1, y_2, \dots, y_m$ . (By assuming limited measurement precision, temperature and pressure can easily be assumed to take on only a finite number of values.) Then we must specify the probabilities  $p_{ij}$  for combinations  $(x_i, y_j)$  for  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, m$ . Hence for temperature and barometric pressure we need the probabilities of all possible combinations.

If we are interested in three random variables, such as outside temperature, barometric pressure, and humidity, we would need probabilities over all possible combinations of the three variables. For more variables, things get progressively more complicated.

There is an important special case where the probability description of several variables simplifies. Two random variables  $x$  and  $y$  are said to be **independent random variables** if the outcome probabilities for one variable do not depend on the outcome of the other. For example, consider the roll of two dice. The probability of an outcome of, say, 4 on the second die is  $1/6$ , no matter what the outcome of the first die. Hence the two random variables corresponding to the spots on the two dice are independent. On the other hand, outside temperature and barometric pressure are not independent, since if pressure is high, temperature is more likely to be high as well.

## Covariance

When considering two or more random variables, their mutual dependence can be summarized conveniently by their **covariance**.

Let  $x_1$  and  $x_2$  be two random variables with expected values  $\bar{x}_1$  and  $\bar{x}_2$ . The covariance of these variables is defined to be

$$\text{cov}(x_1, x_2) = E[(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)].$$

The covariance of two random variables  $x$  and  $y$  is denoted  $\sigma_{xy}$ . Hence for random variables  $x_1$  and  $x_2$  we write  $\text{cov}(x_1, x_2) = \sigma_{x_1, x_2}$  or, alternatively,  $\text{cov}(x_1, x_2) = \sigma_{12}$ . Note that, by symmetry,  $\sigma_{12} = \sigma_{21}$ .

Analogous to (6.2), there is an alternative shorter formula for covariance that is easily derived; namely,

$$\text{cov}(x_1, x_2) = E(x_1 x_2) - \bar{x}_1 \bar{x}_2. \quad (6.3)$$

This is useful in computations.

If two random variables  $x_1$  and  $x_2$  have the property that  $\sigma_{12} = 0$ , then they are said to be **uncorrelated**. This is the situation (roughly) where knowledge of the value of one variable gives no information about the other. If two random variables are independent, then they are uncorrelated. If  $\sigma_{12} > 0$ , the two variables are said to be **positively correlated**. In this case, if one variable is above its mean, the other is likely to be above its mean as well. On the other hand, if  $\sigma_{12} < 0$ , the two variables are said to be **negatively correlated**.

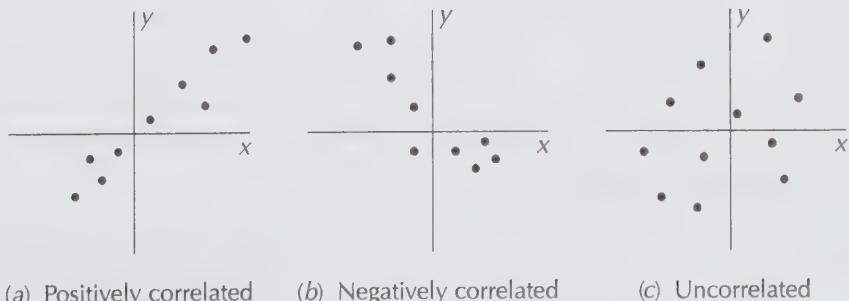
Figure 6.2 illustrates the concept of correlation by showing collections of random samples of two variables  $x$  and  $y$  under the conditions (a) positive correlation, (b) negative correlation, and (c) no correlation.

The following result gives an important bound on the covariance.

**Covariance bound** *The covariance of two random variables satisfies*

$$|\sigma_{12}| \leq \sigma_1 \sigma_2.$$

In the preceding inequality, if  $\sigma_{12} = \sigma_1 \sigma_2$ , the variables are **perfectly correlated**. In this situation, the covariance is as large as possible for the given variances. If one variable were a fixed positive multiple of the other, the two would be perfectly correlated. Conversely, if  $\sigma_{12} = -\sigma_1 \sigma_2$ , the two variables exhibit **perfect negative correlation**.



**FIGURE 6.2 Correlations of data.** Samples are drawn of the pair of random variables  $x$  and  $y$ , and these pairs are plotted on an  $x$ – $y$  diagram. A typical pattern of points obtained is shown in the three cases: (a) positive correlation, (b) negative correlation, and (c) no correlation.

Another useful construct is the **correlation coefficient** of two variables, defined as

$$\rho_{12} = \frac{\sigma_{12}}{\sigma_1 \sigma_2}.$$

From the covariance bound above, we see that  $|\rho_{12}| \leq 1$ .

Note that the variance of a random variable  $x$  is the covariance of that variable with itself. Hence we write  $\sigma_x^2 = \sigma_{xx}$ .

## Variance of a Sum

When we know the covariance between two random variables, it is possible to compute the variance of the sum of the variables. This is a computation that is used frequently in what follows.

Suppose that  $x$  and  $y$  are random variables. We have, by linearity, that  $E(x + y) = \bar{x} + \bar{y}$ . Also by definition,

$$\begin{aligned}\text{var}(x + y) &= E[(x - \bar{x} + y - \bar{y})^2] \\ &= E[(x - \bar{x})^2] + 2E[(x - \bar{x})(y - \bar{y})] + E[(y - \bar{y})^2] \\ &= \sigma_x^2 + 2\sigma_{xy} + \sigma_y^2.\end{aligned}\tag{6.4}$$

This formula is easy to remember because it looks similar to the standard expression for the square of the sum of two algebraic quantities. We just substitute variance for the square and the covariance for the product.

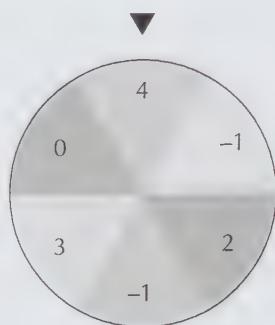
An important special case is where the two variables are uncorrelated. In that case  $\sigma^2 = \sigma_x^2 + \sigma_y^2$ .

**Example 6.4 (Two rolls of the die)** Suppose that a die is rolled twice and the average of the two numbers of spots is recorded as a quantity  $z$ . What are the mean value and the variance of  $z$ ? We let  $x$  and  $y$  denote the values obtained on the first and second rolls, respectively. Then  $z = \frac{1}{2}(x + y)$ . Also  $x$  and  $y$  are uncorrelated, since the rolls of the die are independent. Therefore  $\bar{z} = \frac{1}{2}(\bar{x} + \bar{y}) = 3.5$ , and  $\text{var}(z) = \frac{1}{4}(\sigma_x^2 + \sigma_y^2) = 2.92/2 = 1.46$ . Hence  $\sigma_z = 1.208$ , which is somewhat smaller than the corresponding 1.71 value for a single roll.

## 6.3 Random Returns

When an asset is originally acquired, its rate of return is usually uncertain. Accordingly, we consider the rate of return  $r$  to be a random variable. For analytical purposes we shall, in this chapter, summarize the uncertainty of the rate of return by its expected value (or mean)  $E(r) \equiv \bar{r}$ , by its variance  $E[(r - \bar{r})^2] \equiv \sigma^2$ , and by its covariance with other assets of interest. We can best illustrate how rates of return are represented by considering a few examples.

**FIGURE 6.3 Wheel of fortune.** If you bet \$1 on the wheel, you will receive the amount equal to the value shown in the segment under the marker after the wheel is spun.



**Example 6.5 (Wheel of fortune)** Consider the wheel of fortune shown in Figure 6.3. It is unlike any wheel you are likely to find in an amusement park since its payoffs are quite favorable. If you bet \$1 on the wheel, the payoff you receive is that shown in the segment corresponding to the landing spot. The chance of landing on a given segment is proportional to the area of the segment. For this wheel the probability of each segment is 1/6.

Let us first compute the mean and the variance of the payoff of the wheel. We denote the payoff of segment  $i$  by  $Q_i$ . Therefore the expected payoff is

$$\bar{Q} = \sum_i p_i Q_i = \frac{1}{6}(4 - 1 + 2 - 1 + 3) = 7/6.$$

The variance can be found from the short formula (6.2) to be

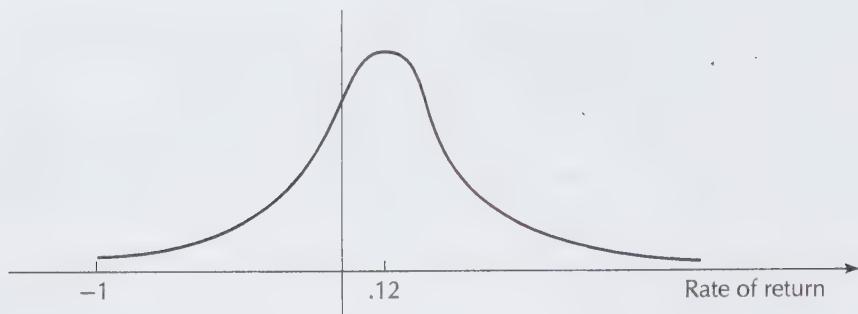
$$\sigma_Q^2 = E(Q^2) - \bar{Q}^2 = \frac{1}{6}(16 + 1 + 4 + 1 + 9) - (7/6)^2 = 3.81.$$

The payoff of the wheel is the same as the total return under the assumption of a \$1 bet. Therefore  $Q = R$  and the rate of return is  $r = Q - 1$ . From this we find

$$\bar{r} = E(r) = \bar{Q} - 1 = 1/6$$

$$\sigma_r^2 = E[(r - \bar{r})^2] = E[(Q - 1 - (\bar{Q} - 1))^2] = \sigma_Q^2 = 3.81.$$

**Example 6.6 (Rate of return on a stock)** Let us consider a share of stock in a major corporation (such as Google, Ford, or Coca-Cola) as an asset. Imagine that we are attempting to describe the rate of return that applies if we were to buy it now and sell it at the end of one year. We ignore transactions costs. As an estimate, we might take  $E(r) = .12$ ; that is, we estimate that the expected rate of return is 12%. This is a reasonable value for the stock of a major corporation, based on the past performance of stocks in the overall market. Now what about the standard deviation? We recognize that the 12% figure is not likely to be hit exactly, and that there can be significant deviations. In fact it is quite possible that the 1-year rate of return could be  $-5\%$  in one year and  $+25\%$  in the next. A reasonable estimate for the standard deviation is about  $.15$ , or 15%. Hence, loosely, we might say that the rate of return is likely to be



**FIGURE 6.4 Probability density of the rate of return of a stock.** The mean rate of return may be about 12% and the standard deviation about 15%. The rate of return cannot be less than  $-1$ .

12% plus or minus 15%. We discuss the process of estimating expected values and standard deviations for stocks in Chapter 9, but this example gives a rough idea of typical magnitudes.

The probability density for the rate of return of this typical stock is shown in Figure 6.4. It has a mean value of  $.12$ , but the return can become arbitrarily large. However, the rate of return can never be less than  $-1$ , since that represents complete loss of the original investment.

**Example 6.7 (Betting wheel)** Two kinds of wheels are useful for the study of investment problems. The wheel of fortune of Example 6.5 is one form of wheel. For that type, one bets on (invests in) the wheel as a whole, and the payoff is determined by the landing segment.

The other kind of wheel is a **betting wheel**, an example of which is shown in Figure 6.5. For this kind of wheel one bets on (invests in) the individual segments of the wheel. For example, for the wheel shown, if one invests \$1 in the white segment, then \$3 will be the payoff if white is the landing segment; otherwise the payoff is zero and the original \$1 is lost. One is allowed to bet different amounts on different segments. A roulette wheel is a betting wheel. From a theoretical viewpoint, a betting wheel is interesting because the returns from different segments are correlated.

For the wheel shown, we may bet on: (1) white, (2) black, or (3) gray, with payoffs 3, 2, or 6, respectively. Note that the bet on white has quite favorable odds.

We can work out the expected rates of return for the three possible bets. It is much easier here to work first with total returns and then subtract 1. For example, for white the return is 3 with probability  $\frac{1}{2}$  and 0 with probability  $\frac{1}{2}$ .

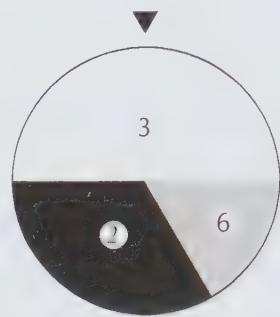
The three expected values are:

$$\bar{R}_1 = \frac{1}{2}(3) + \frac{1}{2}(0) = \frac{3}{2}$$

$$\bar{R}_2 = \frac{1}{3}(2) + \frac{2}{3}(0) = \frac{2}{3}$$

$$\bar{R}_3 = \frac{1}{6}(6) + \frac{5}{6}(0) = 1.$$

**FIGURE 6.5 Betting wheel.** It is possible to bet on any segment of the wheel. If that segment is chosen by the spin, the better receives the amount indicated times the bet.



Likewise, the three variances are, from equation (6.2),

$$\sigma_1^2 = \frac{1}{2}(3^2) - (\frac{3}{2})^2 = 2.25$$

$$\sigma_2^2 = \frac{1}{3}(2^2) - (\frac{2}{3})^2 = .889$$

$$\sigma_3^2 = \frac{1}{6}6^2 - 1 = 5.$$

Finally, we can calculate the covariances using equation (6.3). The expected value of products such as  $E(R_1 R_2)$  are all zero, so we easily find

$$\sigma_{12} = -\frac{3}{2}(\frac{2}{3}) = -1.0$$

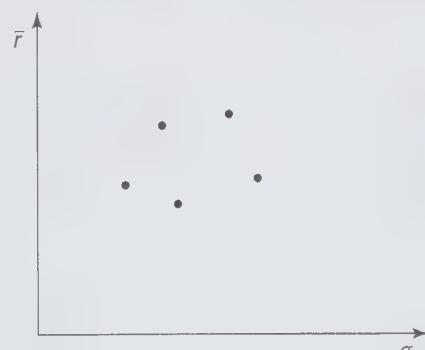
$$\sigma_{13} = -\frac{3}{2}(1) = -1.5$$

$$\sigma_{23} = -\frac{2}{3}(1) = -.67.$$

## Mean–Standard Deviation Diagram

The random rates of return of assets can be represented on a two-dimensional diagram, as shown in Figure 6.6. An asset with mean rate of return  $\bar{r}$  [or  $m$  or  $E(r)$ ] and standard deviation  $\sigma$  is represented as a point in this diagram. The horizontal axis is used for the standard deviation, and the vertical axis is used for

**FIGURE 6.6 Mean–standard deviation diagram.** Assets are described as points on the diagram.



the mean. This diagram is called a mean-standard deviation diagram, or simply  $\bar{r} - \sigma$  diagram.

In such a diagram the standard deviation, rather than the variance, is used as the horizontal axis. This gives both axes comparable units (such as percent per year). Such diagrams are used frequently in mean–variance investment analysis.

## 6.4 Portfolio Mean and Variance

Now that we have the concepts of expected value (or mean) and variance for returns of individual assets and covariances between pairs of assets, we show how these can be used to determine the corresponding mean and variance of the return of a portfolio.

### Mean Return of a Portfolio

Suppose that there are  $n$  assets with (random) rates of return  $r_1, r_2, \dots, r_n$ . These have expected values  $E(r_1) = \bar{r}_1, E(r_2) = \bar{r}_2, \dots, E(r_n) = \bar{r}_n$ .

Suppose that, as in Section 6.1, we form a portfolio of these  $n$  assets using the weights  $w_i, i = 1, 2, \dots, n$ . The rate of return of the portfolio in terms of the return of the individual returns is

$$r = w_1 r_1 + w_2 r_2 + \cdots + w_n r_n.$$

We may take the expected values of both sides, and using linearity (property 2 of the expected value in Section 6.2), we obtain

$$E(r) = w_1 E(r_1) + w_2 E(r_2) + \cdots + w_n E(r_n).$$

In other words, the expected rate of return of the portfolio is found by taking the weighted sum of the individual expected rates of return. So, finding the expected return of a portfolio is easy once we have the expected rates of return of the individual assets from which the portfolio is composed.

### Variance of Portfolio Return

Now let us determine the variance of the rate of return of the portfolio.

We denote the variance of the return of asset  $i$  by  $\sigma_i^2$ , the variance of the return of the portfolio by  $\sigma^2$ , and the covariance of the return of asset  $i$  with asset  $j$  by  $\sigma_{ij}$ . We perform a straightforward calculation:

$$\begin{aligned}\sigma^2 &= E[(r - \bar{r})^2] \\ &= E\left[\left(\sum_{i=1}^n w_i r_i - \sum_{i=1}^n w_i \bar{r}_i\right)^2\right]\end{aligned}$$

$$\begin{aligned}
 &= E \left[ \left( \sum_{i=1}^n w_i (r_i - \bar{r}_i) \right) \left( \sum_{j=1}^n w_j (r_j - \bar{r}_j) \right) \right] \\
 &= E \left[ \sum_{i,j=1}^n w_i w_j (r_i - \bar{r}_i) (r_j - \bar{r}_j) \right] \\
 &= \sum_{i,j=1}^n w_i w_j \sigma_{ij}.
 \end{aligned}$$

This important result shows how the variance of a portfolio's return can be calculated easily from the covariances of the pairs of asset returns and the asset weights used in the portfolio. (Recall,  $\sigma_{ii} = \sigma_i^2$ .)

**Example 6.8 (Two-asset portfolio)** Suppose that there are two assets with  $\bar{r}_1 = .12$ ,  $\bar{r}_2 = .15$ ,  $\sigma_1 = .20$ ,  $\sigma_2 = .18$ , and  $\sigma_{12} = .01$  (values typical for two stocks). A portfolio is formed with weights  $w_1 = .25$  and  $w_2 = .75$ . We can calculate the mean and the variance of the portfolio. First we have the mean,

$$\bar{r} = .25(.12) + .75(.15) = .1425.$$

Second we calculate the variance,

$$\sigma^2 = (.25)^2(.20)^2 + .25(.75)(.01) + .75(.25)(.01) + (.75)^2(.18)^2 = .024475.$$

Note that the two cross terms are equal (since  $w_i w_j = w_j w_i$ ). Hence,

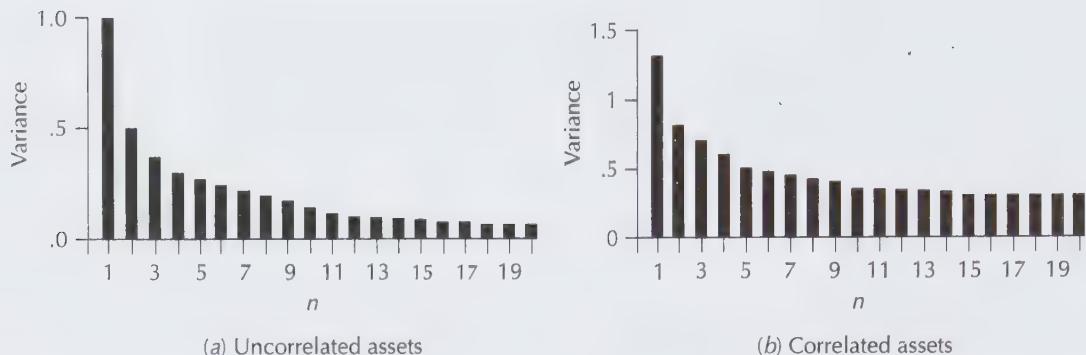
$$\sigma = .1564.$$

## Diversification\*

Portfolios with only a few assets may be subject to a high degree of risk, represented by a relatively large variance. As a general rule, the variance of the return of a portfolio can be reduced by including additional assets in the portfolio, a process referred to as **diversification**. This process reflects the maxim, “Don’t put all your eggs in one basket.”

The effects of diversification can be quantified by using the formulas for combining variances. Suppose as an example that there are many assets, all of which are mutually uncorrelated. That is, the return of each asset is uncorrelated with that of any other asset in the group. Suppose also that the rate of return of each of these assets has mean  $m$  and variance  $\sigma^2$ . Now suppose that a portfolio is constructed by taking equal portions of  $n$  of these assets; that is,  $w_i = 1/n$  for each  $i$ . The overall rate of return of this portfolio is

$$r = \frac{1}{n} \sum_{i=1}^n r_i.$$



**FIGURE 6.7 Effects of diversification.** If assets are uncorrelated, the variance of a portfolio can be made very small. If assets are positively correlated, there is likely to be a lower limit to the variance that can be achieved.

The mean value of this is  $\bar{r} = m$ , which is independent of  $n$ . The corresponding variance is

$$\text{var}(r) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n},$$

where we have used the fact that the individual returns are uncorrelated. The variance decreases rapidly as  $n$  increases, as shown in Figure 6.7(a). This chart shows the variance as a function of  $n$ , the number of assets (when  $\sigma^2 = 1$ ). Note that considerable improvement is obtained by including about six uncorrelated assets.

The situation is somewhat different if the returns of the available assets are correlated. As a simple example suppose again that each asset has a rate of return with mean  $m$  and variance  $\sigma^2$ , but now each return pair has a covariance of  $\text{cov}(r_i, r_j) = .3\sigma^2$  for  $i \neq j$ . Again we form a portfolio by taking equal portions of  $n$  of these assets. In this case,

$$\begin{aligned} \text{var}(r) &= E \left[ \left( \sum_{i=1}^n \frac{1}{n} (r_i - \bar{r}) \right)^2 \right] \\ &= \frac{1}{n^2} E \left\{ \left[ \sum_{i=1}^n (r_i - \bar{r}) \right] \left[ \sum_{j=i}^n (r_j - \bar{r}) \right] \right\} \\ &= \frac{1}{n^2} \sum_{i,j} \sigma_{ij} = \frac{1}{n^2} \left\{ \sum_{i=j} \sigma_{ij} + \sum_{i \neq j} \sigma_{ij} \right\} \\ &= \frac{1}{n^2} \{n\sigma^2 + .3(n^2 - n)\sigma^2\} \end{aligned}$$

$$\begin{aligned}
 &= \frac{\sigma^2}{n} + .3\sigma^2 \left(1 - \frac{1}{n}\right) \\
 &= \frac{.7\sigma^2}{n} + .3\sigma^2.
 \end{aligned}$$

This result is shown in Figure 6.7(b) (where again  $\sigma^2 = 1$ ). In this case it is impossible to reduce the variance below  $.3\sigma^2$ , no matter how large  $n$  is made.

This analysis of diversification is somewhat crude, for we have assumed that all expected rates of return are equal. In general, diversification may reduce the overall expected return while reducing the variance. Most people do not want to sacrifice much expected return for a small decrease in variance, so blind diversification, without an understanding of its influence on both the mean and the variance of return, is not necessarily desirable. This is the motivation behind the general mean–variance approach developed by Markowitz. It makes the trade-offs between mean and variance explicit.

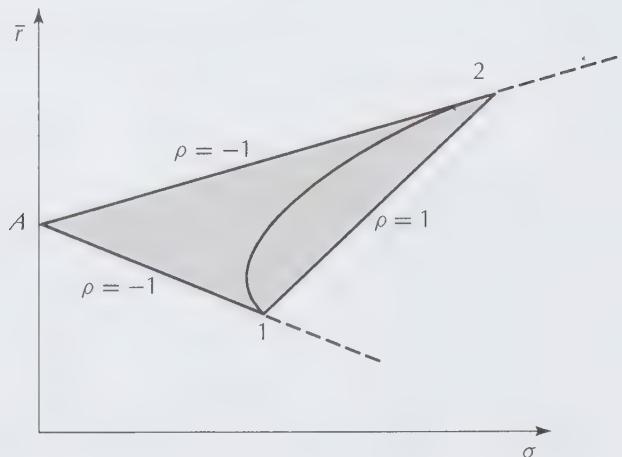
Nevertheless, there is an important lesson to be learned from this simple analysis. Namely, if returns are uncorrelated, it is possible through diversification to reduce portfolio variance essentially to zero by taking  $n$  large. Conversely, if returns are positively correlated, it is more difficult to reduce variance, and there may be a lower limit to what can be achieved.

## Diagram of a Portfolio

Suppose that two assets are represented on a mean–standard deviation diagram. These two assets can be combined, according to some weights, to form a portfolio—a new asset. The mean value and the standard deviation of the rate of return of this new asset can be calculated from the mean, variances, and covariances of the returns of the original assets. However, since covariances are not shown on the diagram, the exact location of the point representing the new asset cannot be determined from the location on the diagram of the original two assets. There are many possibilities, depending on the covariance of these asset returns.

We analyze the possibilities as follows. We begin with two assets as indicated in Figure 6.8. We then define a whole family of portfolios by introducing the variable  $\alpha$ , which defines weights as  $w_1 = 1 - \alpha$  and  $w_2 = \alpha$ . Thus as  $\alpha$  varies from 0 to 1, the portfolio goes from one that contains only asset 1 to one that contains a mixture of assets 1 and 2, and then to one that contains only asset 2. Values of  $\alpha$  outside the range  $0 \leq \alpha \leq 1$  make one or the other of the weights negative, corresponding to short selling.

As  $\alpha$  varies, the new portfolios trace out a curve that includes assets 1 and 2. This curve will look something like the curved shape shown in Figure 6.8, but its exact shape depends on  $\sigma_{12}$ . The solid portion of the curve corresponds to positive combinations of the two assets; the dashed portion corresponds to the shorting of one of them (the one at the opposite end of the solid curve). It can be shown in fact that the solid portion of the curve must lie within the shaded region shown in the figure; that is, it must lie within a triangular region defined by the vertices 1, 2, and a point



**FIGURE 6.8 Combinations of two assets.** When two assets are combined in various combinations, the resulting portfolios sweep out a curve between the points representing the original assets. This curve must lie in the shaded triangular region shown. The edges are defined by the extreme values of the correlation coefficient  $\rho$ , which satisfies  $-1 \leq \rho \leq 1$ .

$A$  on the vertical axis. We state this property formally, but it is not essential that you absorb the details at first reading. It is only necessary to understand the general shape of the curve.

**Portfolio diagram lemma** *The curve in an  $\bar{r} - \sigma$  diagram defined by nonnegative mixtures of two assets 1 and 2 lies within the triangular region defined by the two original assets and the point on the vertical axis of height  $A = (\bar{r}_1\sigma_2 + \bar{r}_2\sigma_1)/(\sigma_1 + \sigma_2)$ .*

**Proof:** The rate of return of the portfolio defined by  $\alpha$  is  $r(\alpha) = (1 - \alpha)r_1 + \alpha r_2$ . The mean value of this return is

$$\bar{r}(\alpha) = (1 - \alpha)\bar{r}_1 + \alpha\bar{r}_2.$$

This says that the mean value is between the original means, in direct proportion to the proportions of the assets. In a 50–50 mix, for example, the new mean will be midway between the original means.

Let us compute the standard deviation of the portfolio. We have, from the general formula of the previous section,

$$\sigma(\alpha) = \sqrt{(1 - \alpha)^2\sigma_1^2 + 2\alpha(1 - \alpha)\sigma_{12} + \alpha^2\sigma_2^2}.$$

Using the definition of the correlation coefficient  $\rho = \sigma_{12}/(\sigma_1\sigma_2)$ , this equation can be written

$$\sigma(\alpha) = \sqrt{(1 - \alpha)^2\sigma_1^2 + 2\rho\alpha(1 - \alpha)\sigma_1\sigma_2 + \alpha^2\sigma_2^2}.$$

This is quite a messy expression. However, we can determine its bounds. We know that  $\rho$  can range over  $-1 \leq \rho \leq 1$ . Using  $\rho = 1$  we find the upper bound

$$\begin{aligned}\sigma(\alpha)^* &= \sqrt{(1-\alpha)^2\sigma_1^2 + 2\alpha(1-\alpha)\sigma_1\sigma_2 + \alpha^2\sigma_2^2} \\ &= \sqrt{[(1-\alpha)\sigma_1 + \alpha\sigma_2]^2} \\ &= (1-\alpha)\sigma_1 + \alpha\sigma_2.\end{aligned}$$

Using  $\rho = -1$  we likewise obtain the lower bound

$$\begin{aligned}\sigma(\alpha)_* &= \sqrt{(1-\alpha)^2\sigma_1^2 - 2\alpha(1-\alpha)\sigma_1\sigma_2 + \alpha^2\sigma_2^2} \\ &= \sqrt{[(1-\alpha)\sigma_1 - \alpha\sigma_2]^2} \\ &= |(1-\alpha)\sigma_1 - \alpha\sigma_2|.\end{aligned}$$

Notice that the upper bound expression is linear in  $\alpha$ , just like the expression for the mean. If we use these two linear expressions, we deduce that both the mean and the standard deviation move proportionally to  $\alpha$  between their values at  $\alpha = 0$  and  $\alpha = 1$ , provided that  $\rho = 1$ . This implies that as  $\alpha$  varies from 0 to 1, the portfolio point will trace out a straight line between the two points. This is the direct line between 1 and 2 indicated in the figure.

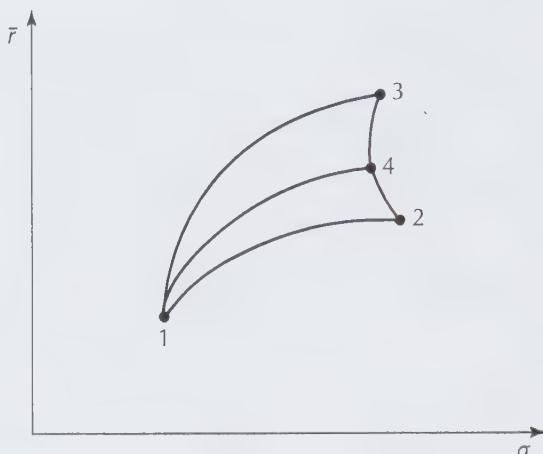
The lower bound expression is nearly linear as well, except for the absolute-value sign. When  $\alpha$  is small, the term inside the absolute-value sign is positive, so we can replace that term by  $(1-\alpha)\sigma_1 - \alpha\sigma_2$ . This remains positive until  $\alpha = \sigma_1/(\sigma_1 + \sigma_2)$ . After that it reverses sign, and so the absolute value becomes  $\alpha\sigma_2 - (1-\alpha)\sigma_1$ . The reversal occurs at the point A given by the expression in the proposition statement. The two linear expressions, together with the linear expression for the mean, imply that the lower bound traces out the kinked line shown in Figure 6.8. We conclude that the curve traced out by the portfolio points must lie within the shaded region; and for an intermediate value of  $\rho$ , it looks like the curve shown. ■

## 6.5 The Feasible Set

Suppose now that there are  $n$  basic assets. We can plot them as points on the mean-standard deviation diagram. Next imagine forming portfolios from these  $n$  assets, using every possible weighting scheme. Hence there are portfolios consisting of each of the  $n$  assets alone, combinations of two assets, combinations of three, and so forth, all the way to arbitrary combinations of all  $n$ . These portfolios are made by letting the weighting coefficients  $w_i$  range over all possible combinations such that  $\sum_{i=1}^n w_i = 1$ .

The set of points that correspond to portfolios is called the **feasible set** or **feasible region**. The feasible set satisfies two important properties.

1. If there are at least three assets (not perfectly correlated and with different means), the feasible set will be a solid two-dimensional region.



**FIGURE 6.9 Three points form a region.** Combinations of assets 2 and 3 sweep out a curve between them. Combination of one of these assets, such as 4, together with asset 1 sweeps out another curve. The family of all these curves forms a solid region.

Figure 6.9 shows why the region will be solid. There are three basic assets: 1, 2, and 3. We know that any two assets define a (curved) line between them as combination portfolios are formed. The three lines between the possible three pairs are shown in Figure 6.9. Now if a combination of, say, assets 2 and 3 is formed to produce asset 4, this can be combined with 1 to form a line connecting 1 and 4. As 4 is moved between 2 and 3, the line between 1 and 4 traces out a solid region.

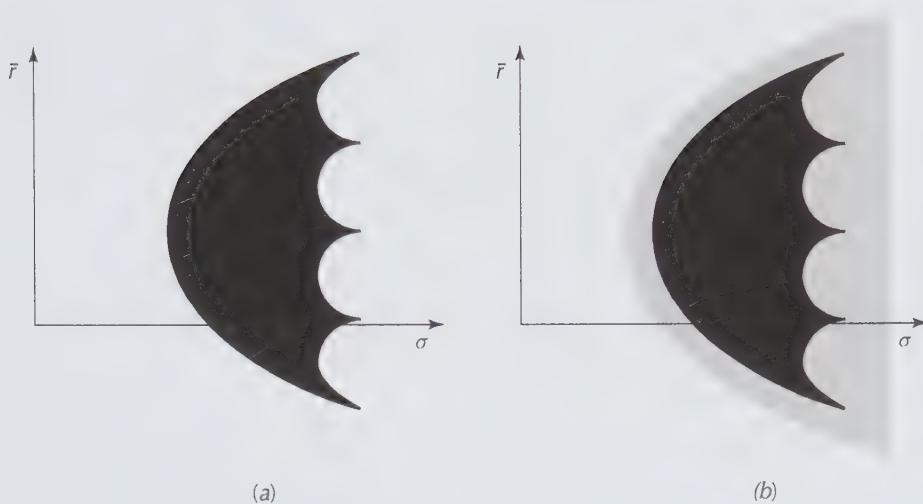
## 2. The feasible region is convex to the left.

This means that given any two points in the region, the straight line connecting them does not cross the left boundary of the feasible set. This follows from the fact that all portfolios (with positive weights) made from two assets lie on or to the left of the line connecting them. A typical feasible region is shown in Figure 6.10(a).

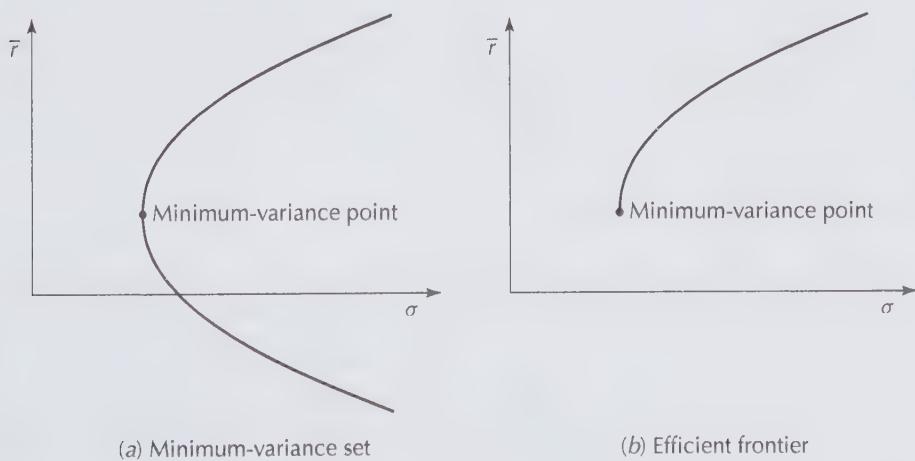
There are two natural, but alternative, definitions of the feasible region, corresponding to whether short selling of assets is allowed or not allowed. The two general conclusions about the shape of the region hold in either case. However, in general the feasible region defined with short selling allowed will contain the region defined without short selling, as shown in Figure 6.10(b). (In general, the leftmost edges of these two regions may partially coincide—unlike the case shown in Figure 6.10.)

## The Minimum-Variance Set and the Efficient Frontier

The left boundary of a feasible set is called the **minimum-variance set**, since for any value of the mean rate of return, the feasible point with the smallest variance (or



**FIGURE 6.10 Feasible region.** The feasible region is the set of all points representing portfolios made from  $n$  original assets. Two such regions can be defined: (a) no shorting and (b) shorting allowed.



**FIGURE 6.11 Special sets.** The minimum-variance set has a characteristic bullet shape. The minimum-variance point is the point with lowest possible variance. The efficient frontier is the upper portion of the minimum-variance set.

standard deviation) is the corresponding left boundary point. The minimum-variance set has a characteristic **bullet** shape, as shown in Figure 6.11 (a). There is a special point on this set having minimum variance. It is termed the **minimum-variance point** (MVP).

Suppose that an investor's choice of portfolio is restricted to the feasible points on a given horizontal line in the  $\bar{r} - \sigma$  plane. All portfolios on this line have the same mean rate of return, but different standard deviations (or variances). Most investors will prefer the portfolio corresponding to the leftmost point on the line; that is, the point with the smallest standard deviation for the given mean. An investor who agrees with this viewpoint is said to be **risk averse**, since he or she seeks to minimize risk (as measured by standard deviation). An investor who would select a point other than the one of minimum standard deviation is said to be **risk preferring**. We direct our analysis to risk-averse investors who, accordingly, prefer to minimize the standard deviation. Such investors are interested in points on the minimum-variance set.

We can turn the argument around 90 degrees and consider portfolios corresponding to the various points on a vertical line; that is, the portfolios with a fixed standard deviation and various mean values. Most investors will prefer the highest point on such a line. In other words, they would select the portfolio of the largest mean for a given level of standard deviation. This property of investors is termed **nonsatiation**, which reflects the idea that, everything else being equal, investors always want more money; hence they want the highest possible expected return for a given standard deviation.

These arguments imply that only the upper part of the minimum-variance set will be of interest to investors who are risk averse and satisfy nonsatiation. This upper portion of the minimum-variance set is termed the **efficient frontier** of the feasible region. It is illustrated in Figure 6.11(b). These are the efficient portfolios, in the sense that they provide the best mean-variance combinations for most investors. We can therefore limit our investigation to this frontier. The next section explains how to calculate points on this frontier.

## 6.6 The Markowitz Model

We are now in a position to formulate a mathematical problem that leads to minimum-variance portfolios. Again assume that there are  $n$  assets. The mean (or expected) rates of return are  $\bar{r}_1, \bar{r}_2, \dots, \bar{r}_n$  and the covariances are  $\sigma_{ij}$ , for  $i, j = 1, 2, \dots, n$ . A portfolio is defined by a set of  $n$  weights  $w_i$ ,  $i = 1, 2, \dots, n$ , that sum to 1. (We allow negative weights, corresponding to short selling.) To find a minimum-variance portfolio, we fix the mean value at some arbitrary value  $\bar{r}$ . Then we find the feasible portfolio of minimum variance that has this mean. Hence we formulate the problem

$$\text{minimize } \frac{1}{2} \sum_{i,j=1}^n w_i w_j \sigma_{ij}$$

$$\begin{aligned} \text{subject to } & \sum_{i=1}^n w_i \bar{r}_i = \bar{r} \\ & \sum_{i=1}^n w_i = 1. \end{aligned}$$

The factor of  $\frac{1}{2}$  in front of the variance is for convenience only. It makes the final form of the equations neater.

The Markowitz problem provides the foundation for single-period investment theory. The problem explicitly addresses the trade-off between expected rate of return and variance of the rate of return in a portfolio. Once the Markowitz problem is formulated, it can be solved numerically to obtain a specific numerical solution. It is also useful to solve the problem analytically because some strong additional conclusions are obtained from the analytic solution. However, as we move to the next chapter, the Markowitz problem is used mainly when a risk-free asset as well as risky assets are available. The existence of a risk-free asset greatly simplifies the nature of the feasible set and also simplifies the analytic solution.

## Solution of the Markowitz Problem\*

We can find the conditions for a solution to this problem using **Lagrange multipliers**  $\lambda$  and  $\mu$ . We form<sup>1</sup> the **Lagrangian**

$$L = \frac{1}{2} \sum_{i,j=1}^n w_i w_j \sigma_{ij} - \lambda \left( \sum_{i=1}^n w_i \bar{r}_i - \bar{r} \right) - \mu \left( \sum_{i=1}^n w_i - 1 \right).$$

We then differentiate the Lagrangian with respect to each variable  $w_i$  and set this derivative to zero.

The differentiation may be a bit difficult if this type of structure is unfamiliar to you. Therefore we shall do it for the two-variable case, after which it will be easy to generalize to  $n$  variables. For two variables,

$$\begin{aligned} L = & \frac{1}{2}(w_1^2 \sigma_1^2 + w_1 w_2 \sigma_{12} + w_2 w_1 \sigma_{21} + w_2^2 \sigma_2^2) \\ & - \lambda(\bar{r}_1 w_1 + \bar{r}_2 w_2 - \bar{r}) - \mu(w_1 + w_2 - 1). \end{aligned}$$

Hence,

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{1}{2}(2\sigma_1^2 w_1 + \sigma_{12} w_2 + \sigma_{21} w_1) - \lambda \bar{r}_1 - \mu \\ \frac{\partial L}{\partial w_2} &= \frac{1}{2}(\sigma_{12} w_1 + \sigma_{21} w_1 + 2\sigma_2^2 w_2) - \lambda \bar{r}_2 - \mu. \end{aligned}$$

Using the fact that  $\sigma_{12} = \sigma_{21}$  and setting these derivatives to zero, we obtain

$$\begin{aligned} \sigma_1^2 w_1 + \sigma_{12} w_2 - \lambda \bar{r}_1 - \mu &= 0 \\ \sigma_{21} w_1 + \sigma_2^2 w_2 - \lambda \bar{r}_2 - \mu &= 0. \end{aligned}$$

---

<sup>1</sup> In general, the Lagrangian is formed by first converting each constraint to one with a zero right-hand side. Then each left-hand side is multiplied by its Lagrange multiplier and subtracted from the objective function. In our problem,  $\lambda$  and  $\mu$  are the multipliers for the first and second constraints, respectively. (See Appendix B.)

This gives us two equations. In addition, there are the two equations of the constraints, so we have a total of four equations. These can be solved<sup>2</sup> for the four unknowns  $w_1$ ,  $w_2$ ,  $\lambda$ , and  $\mu$ .

The general form for  $n$  variables now can be written by obvious generalization. We state the conditions here:

**Equations for efficient set** *The  $n$  portfolio weights  $w_i$  for  $i = 1, 2, \dots, n$  and the two Lagrange multipliers  $\lambda$  and  $\mu$  for an efficient portfolio (with short selling allowed) having mean rate of return  $\bar{r}$  satisfy*

$$\sum_{j=1}^n \sigma_{ij} w_j - \lambda \bar{r}_i - \mu = 0 \quad \text{for } i = 1, 2, \dots, n \quad (6.5a)$$

$$\sum_{i=1}^n w_i \bar{r}_i = \bar{r} \quad (6.5b)$$

$$\sum_{i=1}^n w_i = 1. \quad (6.5c)$$

We have  $n$  equations in (6.5a), plus the two equations of the constraints (6.5b) and (6.5c), for a total of  $n + 2$  equations. Correspondingly, there are  $n + 2$  unknowns: the  $w_i$ 's,  $\lambda$ , and  $\mu$ . The solution to these equations will produce the weights for an efficient portfolio with mean  $\bar{r}$ . Notice that all  $n + 2$  equations are linear, so they can be solved with linear algebra methods.

**Example 6.9 (Three uncorrelated assets)** Suppose there are three uncorrelated assets. Each has variance 1, and the mean values are 0.1, 0.2, and 0.3, respectively. There is a bit of simplicity and symmetry in this situation, which makes it relatively easy to find an explicit solution.

We have  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = 1$  and  $\sigma_{12} = \sigma_{23} = \sigma_{13} = 0$ . For convenience we define  $\lambda$  to be 0.1  $\bar{\lambda}$ , where  $\bar{\lambda}$  is the actual Lagrange multiplier. Then equations (6.5a–c) become

$$\begin{aligned} w_1 - \lambda - \mu &= 0 \\ w_2 - 2\lambda - \mu &= 0 \\ w_3 - 3\lambda - \mu &= 0 \\ w_1 + 2w_2 + 3w_3 &= 10\bar{r} \\ w_1 + w_2 + w_3 &= 1. \end{aligned}$$

---

<sup>2</sup> The case of two assets is actually degenerate because the two unknowns  $w_1$  and  $w_2$  are uniquely determined by the two constraints. The degeneracy (usually) disappears when there are three or more assets. Nevertheless, the equations obtained for the two-asset case foreshadow the pattern of the corresponding equations for  $n$  assets.

The top three equations can be solved for  $w_1$ ,  $w_2$ , and  $w_3$ , in terms of  $\lambda$  and  $\mu$ , and substituted into the bottom two equations. This leads to

$$14\lambda + 6\mu = 10\bar{r}$$

$$6\lambda + 3\mu = 1.$$

These two equations can be solved to yield  $\lambda = 5\bar{r} - 1$  and  $\mu = 2\frac{1}{3} - 10\bar{r}$ . Then

$$w_1 = \frac{4}{3} - 5\bar{r}$$

$$w_2 = \frac{1}{3}$$

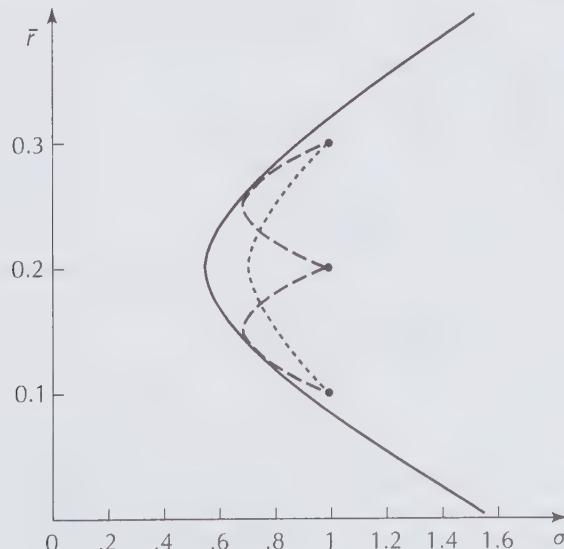
$$w_3 = 5\bar{r} - \frac{2}{3}.$$

The standard deviation at the solution is  $\sqrt{w_1^2 + w_2^2 + w_3^2}$ , which by direct substitution gives

$$\sigma = \sqrt{\frac{7}{3} - 20\bar{r} + 5\bar{r}^2}. \quad (6.6)$$

The minimum-variance point is, by symmetry, at  $\bar{r} = 0.2$ , with  $\sigma = \sqrt{3}/3 = 0.58$ . The feasible region is the region bounded by the bullet-shaped curve shown in Figure 6.12.

The foregoing analysis assumes that shorting of assets is allowed. If shorting is not allowed, the feasible set will be smaller, as discussed in the next subsection.



**FIGURE 6.12 Three-asset example.** The feasible region with shorting contains the feasible region without shorting. The outside curve is the minimum-variance set with shorting allowed. The short curved lines are portfolios made up of two of the assets at a time.

## Nonnegativity Constraints\*

In the preceding derivation, the signs of the  $w_i$  variables were not restricted, which meant that short selling was allowed. We can prohibit short selling by restricting each  $w_i$  to be nonnegative. This leads to the following alternative statement of the Markowitz problem:

$$\text{minimize} \quad \frac{1}{2} \sum_{i,j=1}^n \sigma_{ij} w_i w_j \quad (6.7a)$$

$$\text{subject to} \quad \sum_{i=1}^n \bar{r}_i w_i = \bar{r} \quad (6.7b)$$

$$\sum_{i=1}^n w_i = 1 \quad (6.7c)$$

$$w_i \geq 0 \quad \text{for } i = 1, 2, \dots, n. \quad (6.7d)$$

This problem cannot be reduced to the solution of a set of linear equations. It is termed a **quadratic program**, since the objective is quadratic and the constraints are linear equalities and inequalities. Special computer programs are available for solving such problems, but small to moderate-sized problems of this type can be solved readily with spreadsheet programs. In the financial industry there are a multitude of special-purpose programs designed to solve this problem for hundreds or even thousands of assets.

A significant difference between the two formulations is that when short selling is allowed, most, if not all, of the optimal  $w_i$ 's have nonzero values (either positive or negative), so essentially all assets are used. By contrast, when short selling is not allowed, typically many weights are equal to zero.

## 6.7 The Two-Fund Theorem\*

The minimum-variance set has an important property that greatly simplifies its computation. Recall that points in this set satisfy the system of  $n+2$  linear equations [eqs. (6.5a – c)], which is repeated here:

$$\sum_{j=1}^n \sigma_{ij} w_j - \lambda \bar{r}_i - \mu = 0 \quad \text{for } i = 1, 2, \dots, n \quad (6.8a)$$

$$\sum_{i=1}^n w_i \bar{r}_i = \bar{r} \quad (6.8b)$$

$$\sum_{i=1}^n w_i = 1. \quad (6.8c)$$

Suppose that there are two known solutions,  $\mathbf{w}^1 = (w_1^1, w_2^1, \dots, w_n^1)$ ,  $\lambda^1$ ,  $\mu^1$  and  $\mathbf{w}^2 = (w_1^2, w_2^2, \dots, w_n^2)$ ,  $\lambda^2$ ,  $\mu^2$ , with expected rates of return  $\bar{r}^1$  and  $\bar{r}^2$ , respectively. Let

us form a combination by multiplying the first by  $\alpha$  and the second by  $(1 - \alpha)$ . By direct substitution, we see that the result is also a solution to the  $n + 2$  equations, corresponding to the expected value  $\alpha\bar{r}^1 + (1 - \alpha)\bar{r}^2$ . To check this in detail, notice that  $\alpha\mathbf{w}^1 + (1 - \alpha)\mathbf{w}^2$  is a legitimate portfolio with weights that sum to 1; hence (6.8c) is satisfied. Next notice that the expected return is in fact  $\alpha\bar{r}_1 + (1 - \alpha)\bar{r}_2$ ; hence (6.8b) is satisfied for that value. Finally, notice that since both solutions make the left side of (6.8a) equal to zero, their combination does also; hence (6.8a) is satisfied. This implies that the combination portfolio  $\alpha\mathbf{w}^1 + (1 - \alpha)\mathbf{w}^2$  is also a solution; that is, it also represents a point in the minimum-variance set. This simple result is usually quite surprising to most people on their first exposure to the subject, but it highlights an important property of the minimum-variance set.

To use this result, suppose  $\mathbf{w}^1$  and  $\mathbf{w}^2$  are two different portfolios in the minimum-variance set. Then as  $\alpha$  varies over  $-\infty < \alpha < \infty$ , the portfolios defined by  $\alpha\mathbf{w}^1 + (1 - \alpha)\mathbf{w}^2$  sweep out the entire minimum-variance set. We can, of course, select the two original solutions to be efficient (that is, on the upper portion of the minimum-variance set), and these will generate all other efficient points (as well as all other points in the minimum-variance set). This result is often stated in a form that has operational significance for investors:

**The two-fund theorem** *Two efficient funds (portfolios) can be established so that any efficient portfolio can be duplicated, in terms of mean and variance, as a combination of these two. In other words, all investors seeking efficient portfolios need only invest in combinations of these two funds.*

This result has dramatic implications. According to the two-fund theorem, two **mutual funds**<sup>3</sup> could provide a complete investment service for everyone. There would be no need for anyone to purchase individual stocks separately; they could just purchase shares in the mutual funds. This conclusion, however, is based on the assumption that everyone cares only about mean and variance; that everyone has the same assessment of the means, variances, and covariances; and that a single-period framework is appropriate. All of these assumptions are quite tenuous. Nevertheless, if you are an investor without the time or inclination to make careful assessments, you might choose to find two funds managed by people whose assessments you trust, and invest in those two funds.

The two-fund theorem also has implications for computation. In order to solve equations (6.5a–c) for all values of  $\bar{r}$  it is only necessary to find two solutions and then form combinations of those two. A particularly simple way to specify two solutions is to specify values of  $\lambda$  and  $\mu$ . Convenient choices are (a)  $\lambda = 0, \mu = 1$  and (b)  $\lambda = 1, \mu = 0$ . In either of these solutions the constraint  $\sum_{i=1}^n w_i = 1$  may be violated, but this can be remedied later by normalizing all  $w_i$ 's by a common scale factor. The solution obtained by choice (a) ignores the constraint on the expected mean rate of return; hence this is the minimum-variance point. The overall procedure is illustrated in the following example.

---

<sup>3</sup> A mutual fund is an investment company that accepts investment capital from individuals and reinvests that capital in a diversity of individual stocks. Each individual is entitled to his or her proportionate share of the funds portfolio value, less certain operating fees and commissions.

**TABLE 6.2**  
A SECURITIES PORTFOLIO

Security	Covariance V					$\bar{r}$
1	2.30	.93	.62	.74	-.23	15.1
2	.93	1.40	.22	.56	.26	12.5
3	.62	.22	1.80	.78	-.27	14.7
4	.74	.56	.78	3.40	-.56	9.02
5	-.23	.26	-.27	-.56	2.60	17.68
	$v^1$	$v^2$	$w^1$	$w^2$		
	.141	3.652	.088	.158		
	.401	3.583	.251	.155		
	.452	7.248	.282	.314		
	.166	.874	.104	.038		
	.440	7.706	.275	.334		
Mean			14.413	15.202		
Variance			.625	.659		
Std. dev.			.791	.812		

The covariances and mean rates of return are shown for five securities. The portfolio  $w^1$  is the minimum-variance point, and  $w^2$  is another efficient portfolio made from these five securities.

**Example 6.10 (A securities portfolio)** The information concerning the 1-year covariances and mean values of the rates of return on five securities is shown in the top part of Table 6.2. The mean values are expressed on a percentage basis, whereas the covariances are expressed in units of  $(\text{percent})^2/100$ . For example, the first security has an expected rate of return of  $15.1\% = .151$  and a variance of return of  $.023$ , which translates into a standard deviation of  $\sqrt{.023} = .152 = 15.2\%$  per year.

We shall find two funds in the minimum-variance set. First we set  $\lambda = 0$  and  $\mu = 1$  in (6.5). We thus solve the system of equations

$$\sum_{j=1}^5 \sigma_{ij} v_j^1 = 1$$

for the vector  $\mathbf{v}^1 = (v_1^1, v_2^1, \dots, v_5^1)$ . This solution can be found using a spreadsheet package that solves linear equations. The coefficients of the equation are those of the covariance matrix, and the right-hand sides are all 1's. The resulting  $v_j^1$ 's are listed in the first column of the bottom part of Table 6.2 as components of the vector  $\mathbf{v}^1$ .

Next we normalize the  $v_i^1$ 's so that they sum to 1, obtaining  $w_i^1$ 's as

$$w_i^1 = \frac{v_i^1}{\sum_{j=1}^5 v_j^1}.$$

The vector  $\mathbf{w}^1 = (w_1^1, w_2^1, \dots, w_5^1)$  defines the minimum-variance point.

Second we set  $\mu = 0$  and  $\lambda = 1$ . We thus solve the system of equations

$$\sum_{j=1}^5 \sigma_{ij} v_j^2 = \bar{r}_i, \quad i = 1, 2, \dots, 5$$

for a solution  $\mathbf{v}^2 = (v_1^2, v_2^2, \dots, v_5^2)$ . Again we normalize the resulting vector  $\mathbf{v}^2$  so its components sum to 1, to obtain  $\mathbf{w}^2$ . The vectors  $\mathbf{v}^1, \mathbf{v}^2, \mathbf{w}^1, \mathbf{w}^2$  are shown in the bottom part of Table 6.2. Also shown are the means, variances, and standard deviations corresponding to the portfolios defined by  $\mathbf{w}^1$  and  $\mathbf{w}^2$ . All efficient portfolios are combinations of these two.

## 6.8 Inclusion of a Risk-Free Asset

In the previous few sections we have implicitly assumed that the  $n$  assets available are all risky; that is, they each have  $\sigma > 0$ . A **risk-free asset** has a return that is deterministic (that is, known with certainty) and therefore has  $\sigma = 0$ . In other words, a risk-free asset is a pure interest-bearing instrument; its inclusion in a portfolio corresponds to lending or borrowing cash at the risk-free rate. Lending (such as the purchase of a bond) corresponds to the risk-free asset having a positive weight, whereas borrowing corresponds to its having a negative weight.

The inclusion of a risk-free asset in the list of possible assets is necessary to obtain realism. Investors invariably have the opportunity to borrow or lend. Fortunately, as we shall see shortly, inclusion of a risk-free asset introduces a mathematical degeneracy that greatly simplifies the shape of the efficient frontier.

To explain the degeneracy condition, suppose that there is a risk-free asset with a (deterministic) rate of return  $r_f$ . Consider any other risky asset with rate of return  $r$ , having mean  $\bar{r}$  and variance  $\sigma^2$ . Note that the covariance of these two returns must be zero. This is because the covariance is defined to be  $E[(r - \bar{r})(r_f - r_f)] = 0$ .

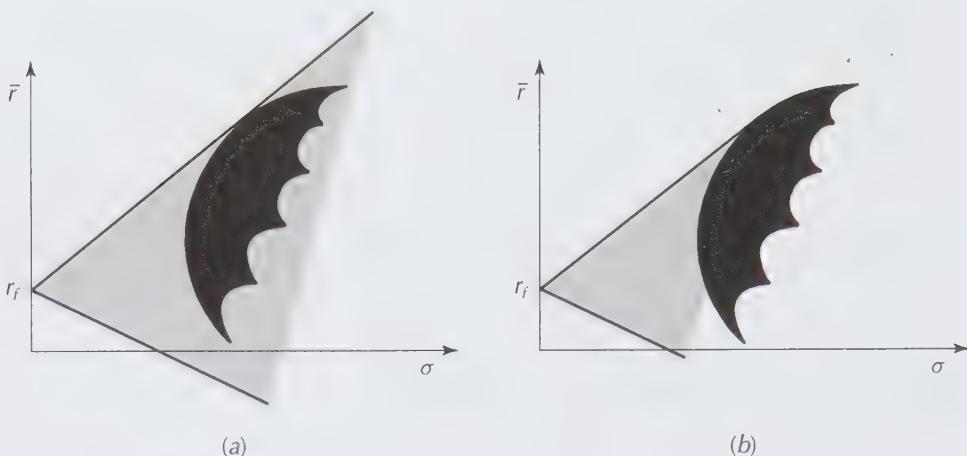
Now suppose that these two assets are combined to form a portfolio using a weight of  $\alpha$  for the risk-free asset and  $1 - \alpha$  for the risky asset, with  $\alpha \leq 1$ . The mean rate of return of this portfolio will be  $\alpha r_f + (1 - \alpha)\bar{r}$ . The standard deviation of the return will be  $\sqrt{(1 - \alpha)^2 \sigma^2} = (1 - \alpha)\sigma$ . This is because the risk-free asset has no variance and no covariance with the risky asset. The only term left in the formula is that due to the risky asset.

Overall, we see that the portfolio rate of return has

$$\begin{aligned} \text{mean} &= \alpha r_f + (1 - \alpha)\bar{r} \\ \text{standard deviation} &= (1 - \alpha)\sigma. \end{aligned}$$

These equations show that both the mean and the standard deviation of the portfolio vary linearly with  $\alpha$ . This means that as  $\alpha$  varies, the point representing the portfolio traces out a straight line in the  $\bar{r}-\sigma$  plane.

Suppose now that there are  $n$  risky assets with known mean rates of return  $\bar{r}_i$  and known covariances  $\sigma_{ij}$ . In addition, there is a risk-free asset with rate of



**FIGURE 6.13 Effect of a risk-free asset.** Inclusion of a risk-free asset adds lines to the feasible region. (a) If both borrowing and lending are allowed, a complete infinite triangular region is obtained. (b) If only lending is allowed, the region will have a triangular front end, but will curve for larger  $\sigma$ .

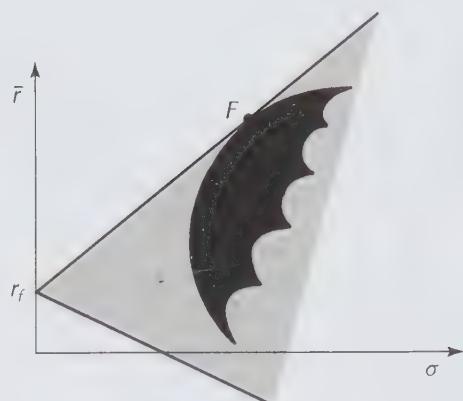
return  $r_f$ . The inclusion of the risk-free asset in the list of available assets has a profound effect on the shape of the feasible region. The reason for this is shown in Figure 6.13(a). First we construct the ordinary feasible region, defined by the  $n$  risky assets. (This region may be either the one constructed with shorting allowed or the one constructed without shorting.) This region is shown as the darkly shaded region in the figure. Next, for each asset (or portfolio) in this region we form combinations with the risk-free asset. In forming these combinations we allow borrowing or lending of the risk-free asset, but only purchase of the risky asset. These new combinations trace out the infinite straight line originating at the risk-free point, passing through the risky asset, and continuing indefinitely. There is a line of this type for every asset in the original feasible set. The totality of these lines forms a triangularly shaped feasible region, indicated by the light shading in the figure.

This is a beautiful result. The feasible region is an infinite triangle whenever a risk-free asset is included in the universe of available assets.

If borrowing of the risk-free asset is not allowed (no shorting of this asset), we can adjoin only the finite line segments between the risk-free asset and points in the original feasible region. We cannot extend these lines further, since this would entail borrowing of the risk-free asset. The inclusion of these finite line segments leads to a new feasible region with a straight-line front edge but a rounded top, as shown in Figure 6.13(b).

**FIGURE 6.14 One-fund theorem.**

When both borrowing and lending at the risk-free rate are allowed, there is a unique fund  $F$  of risky assets that is efficient. All points on the efficient frontier are combinations of  $F$  and the risk-free asset.



## 6.9 The One-Fund Theorem

When risk-free borrowing and lending are available, the efficient set consists of a single straight line, which is the top of the triangular feasible region. This line is tangent to the original feasible set of risky assets. (See Figure 6.14.) There will be a point  $F$  in the original feasible set that is on the line segment defining the overall efficient set. It is clear that *any* efficient point (any point on the line) can be expressed as a combination of this asset and the risk-free asset. We obtain different efficient points by changing the weighting between these two (including negative weights of the risk-free asset to borrow money in order to leverage the buying of the risky asset). The portfolio represented by the tangent point can be thought of as a fund made up of assets and sold as a unit. The role of this fund is summarized by the following statement:

**The one-fund theorem** *There is a single fund  $F$  of risky assets such that any efficient portfolio can be constructed as a combination of the fund  $F$  and the risk-free asset.*

This is a final conclusion of mean–variance portfolio theory, and this conclusion is the launch point for the next chapter. It is fine to stop reading here, and (after doing some exercises) go on to the next chapter. But if you want to see how to calculate the special efficient point  $F$ , read the specialized subsection that follows.

### Solution Method\*

How can we find the tangent point that represents the efficient fund? We just characterize that point in terms of an optimization problem. Given a point in the feasible region, we draw a line between the risk-free asset and that point. We denote the angle between that line and the horizontal axis by  $\theta$ . For any feasible (risky) portfolio  $p$ , we have

$$\tan \theta = \frac{\bar{r}_p - r_f}{\sigma_p}.$$

The tangent portfolio is the feasible point that maximizes  $\theta$  or, equivalently, maximizes  $\tan \theta$ . It turns out that this problem can be reduced to the solution of a system of linear equations.

To develop the solution, suppose, as usual, that there are  $n$  risky assets. We assign weights  $w_1, w_2, \dots, w_n$  to the risky assets such that  $\sum_{i=1}^n w_i = 1$ . There is zero weight on the risk-free asset in the tangent fund. (Note that we are allowing short selling among the risky assets.) For  $r_p = \sum_{i=1}^n w_i r_i$ , we have  $\bar{r}_p = \sum_{i=1}^n w_i \bar{r}_i$  and  $r_f = \sum_{i=1}^n w_i r_f$ . Thus,

$$\tan \theta = \frac{\sum_{i=1}^n w_i (\bar{r}_i - r_f)}{\left( \sum_{i,j=1}^n \sigma_{ij} w_i w_j \right)^{1/2}}.$$

It should be clear that multiplication of all  $w_i$ 's by a constant will not change the expression, since the constant will cancel. Hence it is not necessary to impose the constraint  $\sum_{i=1}^n w_i = 1$  here.

We then set the derivative of  $\tan \theta$  with respect to each  $w_k$  equal to zero. This leads (see Exercise 10) to the following equations:

$$\sum_{i=1}^n \sigma_{ki} \lambda w_i = \bar{r}_k - r_f, \quad k = 1, 2, \dots, n, \quad (6.9)$$

where  $\lambda$  is an (unknown) constant. Making the substitution  $v_i = \lambda w_i$  for each  $i$ , (6.9) becomes

$$\sum_{i=1}^n \sigma_{ki} v_i = \bar{r}_k - r_f, \quad k = 1, 2, \dots, n. \quad (6.10)$$

We solve these linear equations for the  $v_i$ 's and then normalize to determine the  $w_i$ 's; that is,

$$w_i = \frac{v_i}{\sum_{k=1}^n v_k}.$$

**Example 6.11 (Three uncorrelated assets)** We consider again Example 6.9, where the three risky assets were uncorrelated and each had variance equal to 1. The three mean rates of return were  $\bar{r}_1 = 0.1$ ,  $\bar{r}_2 = 0.2$ , and  $\bar{r}_3 = 0.3$ . We assume in addition that there is a risk-free asset with rate  $r_f = 0.05$ .

We apply (6.9), which is very simple in this case because the covariances are all zero, to find

$$v_1 = .1 - .05 = .05$$

$$v_2 = .2 - .05 = .15$$

$$v_3 = .3 - .05 = .25.$$

We then normalize these values by dividing by their sum, .45, and find

$$w_1 = \frac{1}{9}, \quad w_2 = \frac{1}{3}, \quad w_3 = \frac{5}{9}.$$

**Example 6.12 (A larger portfolio)** Consider the five risky assets of Example 6.10. Assume also that there is a risk-free asset with  $r_f = 10\%$ . We can easily find the special fund  $F$  by using the fact that portfolio  $F$  is a combination of two known efficient points.

We note that the system of equations (6.10) is identical to those used to find  $\mathbf{v}^1$  and  $\mathbf{v}^2$  in Example 6.10, but with a different right-hand side. In fact, the solution to equation (6.10) is  $\mathbf{v} = \mathbf{v}^2 - r_f \mathbf{v}^1$ . Thus (using  $r_f = 10$  to be consistent with the units in the earlier example),  $\mathbf{v} = (2.242, -0.427, 2.728, -0.786, 3.306)$ . We normalize this to obtain the final result:  $\mathbf{w} = (.317, -0.060, .386, -.111, .468)$ .

There are abnormal cases where there is no Markowitz portfolio that can serve as the One-Fund. These situations are signaled by the fact that the  $v_i$ 's cannot be normalized, because their sum is negative or, worse yet, zero.

**Example 6.13** Suppose there are two risky assets that are uncorrelated and that each has variance 1. The corresponding expected rates of return are  $\bar{r}_1 = .13$  and  $\bar{r}_2 = .09$ . The risk-free asset has rate of return  $r_f = .11$ . In this case

$$v_1 = .13 - .11 = .02$$

$$v_2 = .09 - .11 = -.02.$$

Hence,  $v_1 + v_2 = 0$ , and these cannot be converted to  $w_1, w_2$  by normalization. Generally, such problems occur if the risk-free rate is greater than the expected rate of return of some of the risky assets.

## Explicit Solution

For the case of two risky assets and one risk-free asset, there are simple explicit formulas for the two optimal weights of the risky assets:

$$w_1 = \frac{[\bar{r}_1 - r_f]\sigma_2^2 - [\bar{r}_2 - r_f]\text{cov}(r_1, r_2)}{[\bar{r}_1 - r_f]\sigma_2^2 + [\bar{r}_2 - r_f]\sigma_1^2 - [\bar{r}_1 - r_f + \bar{r}_2 - r_f]\text{cov}(r_1, r_2)} \quad (6.11)$$

$$w_2 = 1 - w_1.$$

## 6.10 Summary

The study of one-period investment situations is based on asset and portfolio returns. Both total returns and rates of return are used. The return of an asset may be uncertain, in which case it is useful to consider it formally as a random variable. The probabilistic properties of such random returns can be summarized by their expected values, their variances, and their covariances with each other.

A portfolio is defined by allocating fractions of initial wealth to individual assets. The fractions (or weights) must sum to 1; but some of these weights may be negative if short selling is allowed. The return of a portfolio is the weighted sum of the returns of its individual assets, with the weights being those that define the portfolio.

The expected return of the portfolio is, likewise, equal to the weighted average of the expected returns of the individual assets. The variance of the portfolio is determined by a more complicated formula:  $\sigma^2 = \sum_{i,j=1}^n w_i w_j \sigma_{ij}$ , where the  $w_i$ 's are the weights and the  $\sigma_{ij}$ 's are the covariances.

From a given collection of  $n$  risky assets, there results a set of possible portfolios made from all possible weights of the  $n$  individual assets. If the mean and the standard deviation of these portfolios are plotted on a diagram with vertical axis  $\bar{r}$  (the mean) and horizontal axis  $\sigma$  (the standard deviation), the region so obtained is called the feasible region. Two alternative feasible regions are defined: one allowing shorting of assets and one not allowing shorting.

It can be argued that investors who measure the value of a portfolio in terms of its mean and its standard deviation, who are risk averse, and who have the nonsatiation property will select portfolios on the upper left-hand portion of the feasible region—the efficient frontier.

Points on the efficient frontier can be characterized by an optimization problem originally formulated by Markowitz. This problem seeks the portfolio weights that minimize variance for a given value of mean return. Mathematically, this is a problem with a quadratic objective and two linear constraints. If shorting is allowed (so that the weights may be negative as well as positive), the optimal weights can be found by solving a system of  $n + 2$  linear equations and  $n + 2$  unknowns. Otherwise if shorting is not allowed, the Markowitz problem can be solved by special quadratic programming packages.

An important property of the Markowitz problem, when shorting is allowed, is that if two solutions are known, then any weighted combination of these two solutions is also a solution. This leads to the fundamental two-fund theorem: investors seeking efficient portfolios need only invest in two master efficient funds.

Usually it is appropriate to assume that, in addition to  $n$  risky assets, there is available a risk-free asset with fixed rate of return  $r_f$ . The inclusion of such an asset greatly simplifies the shape of the feasible region, transforming the upper boundary into a straight line. This line is the efficient frontier. The straight-line frontier touches the original feasible region (the region defined by the risky assets only) at a single point  $F$ . This leads to the important one-fund theorem: investors seeking efficient portfolios need only invest in one master fund of risky assets and in the risk-free asset. Different investors may prefer different combinations of these two.

The single efficient fund of risky assets  $F$  can be found by solving a system of  $n$  linear equations and  $n$  unknowns. When the solution to this system is normalized so that its components sum to 1, the resulting components are the weights of the risky assets in the master fund.

## Exercises

1. (Shorting with margin) Suppose that to short a stock you are required to deposit an amount equal to the initial price  $X_0$  of the stock. At the end of 1 year the stock price is  $X_1$  and you liquidate your position. You receive your profit from shorting equal to  $X_0 - X_1$  and you recover your original deposit. If  $R$  is the total return of the stock, what is the total return on your short?

2. (Dice product) Two dice are rolled and the two resulting values are multiplied together to form the quantity  $z$ . What are the expected value and the variance of the random variable  $z$ ? [Hint: Use the independence of the two separate dice.]

**TABLE 6.3**  
**TWO CORRELATED CASES**

Asset	$\bar{r}$	$\sigma$
A	10.0%	15%
B	18.0%	30%

3. (Two correlated assets) The correlation  $\rho$  between assets A and B is .1, and other data are given in Table 6.3. [Note:  $\rho = \sigma_{AB}/(\sigma_A \sigma_B)$ .]
- Find the proportions  $\alpha$  of A and  $(1 - \alpha)$  of B that define a portfolio of A and B having minimum standard deviation.
  - What is the value of this minimum standard deviation?
  - What is the expected return of this portfolio?
4. (Two stocks) Two stocks are available. The corresponding expected rates of return are  $\bar{r}_1$  and  $\bar{r}_2$ ; the corresponding variances and covariances are  $\sigma_1^2$ ,  $\sigma_2^2$ , and  $\sigma_{12}$ . What percentages of total investment should be invested in each of the two stocks to minimize the total variance of the rate of return of the resulting portfolio? What is the mean rate of return of this portfolio?
5. (Rain insurance) Gavin Jones's friend is planning to invest \$1 million in a rock concert to be held 1 year from now. The friend figures that he will obtain \$3 million revenue from his \$1 million investment—unless, my goodness, it rains. If it rains, he will lose his entire investment. There is a 50% chance that it will rain the day of the concert. Gavin suggests that he buy rain insurance. He can buy one unit of insurance for \$.50, and this unit pays \$1 if it rains and nothing if it does not. He may purchase as many units as he wishes, up to \$3 million.
- What is the expected rate of return on his investment if he buys  $u$  units of insurance? (The cost of insurance is in addition to his \$1 million investment.)
  - What number of units will minimize the variance of his return? What is this minimum value? And what is the corresponding expected rate of return? [Hint: Before calculating a general expression for variance, think about a simple answer.]
6. (Wild cats) Suppose there are  $n$  assets which are uncorrelated. (They might be  $n$  different “wild cat” oil well prospects.) You may invest in any one, or in any combination of them. The mean rate of return  $\bar{r}$  is the same for each asset, but the variances are different. The return on asset  $i$  has a variance of  $\sigma_i^2$  for  $i = 1, 2, \dots, n$ .
- Show the situation on an  $\bar{r} - \sigma$  diagram. Describe the efficient set.
  - Find the minimum-variance point. Express your result in terms of

$$\bar{\sigma}^2 = \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} \right)^{-1}.$$

7. (Markowitz fun) There are just three assets with rates of return  $r_1, r_2$ , and  $r_3$ , respectively. The covariance matrix and the expected rates of return are

$$\mathbf{V} = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}, \quad \bar{\mathbf{r}} = \begin{bmatrix} 4 \\ 8 \\ 8 \end{bmatrix}.$$

- (a) Find the minimum-variance portfolio. [Hint: By symmetry  $w_1 = w_3$ .]
- (b) Find another efficient portfolio by setting  $\lambda = 1, \mu = 0$ .
- (c) If the risk-free rate is  $r_f = .2$ , find the efficient portfolio of risky assets.

8. (Tracking) Suppose that it is impractical to use all the assets that are incorporated into a specified portfolio (such as a given efficient portfolio). One alternative is to find the portfolio, made up of a given set of  $n$  stocks, that tracks the specified portfolio most closely—in the sense of minimizing the variance of the difference in returns.

Specifically, suppose that the target portfolio has (random) rate of return  $r_M$ . Suppose that there are  $n$  assets with (random) rates of return  $r_1, r_2, \dots, r_n$ . We wish to find the portfolio rate of return

$$r = \alpha_1 r_1 + \alpha_2 r_2 + \cdots + \alpha_n r_n$$

(with  $\sum_{i=1}^n \alpha_i = 1$ ) minimizing  $\text{var}(r - r_M)$ .

- (a) Find a set of equations for the  $\alpha_i$ 's.
- (b) Although this portfolio tracks the desired portfolio most closely in terms of variance, it may sacrifice the mean. Hence a logical approach is to minimize the variance of the tracking error subject to achieving a given mean return. As the mean is varied, this results in a family of portfolios that are efficient in a new sense—say, tracking efficient. Find the equation for the  $\alpha_i$ 's that are tracking efficient.

9. (Betting wheel) Consider a general betting wheel with  $n$  segments. The payoff for a \$1 bet on a segment  $i$  is  $A_i$ . Suppose you bet an amount  $B_i = 1/A_i$  on segment  $i$  for each  $i$ . Show that the amount you win is independent of the outcome of the wheel. What is the risk-free rate of return for the wheel? Apply this to the wheel in Example 6.7.

10. (Efficient portfolio  $\diamond$ ) Derive equation (6.9). [Hint: Note that

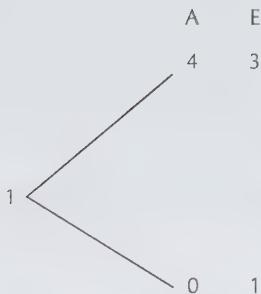
$$\frac{\partial}{\partial w_i} \left( \sum_{ij}^n \sigma_{ij} w_i w_j \right)^{1/2} = \left( \sum_{ij}^n \sigma_{ij} w_i w_j \right)^{-1/2} \sum_{j=1}^n \sigma_{ij} w_j.$$

11. (Two similar assets) Two assets with expected rates of return  $\bar{r}_1$  and  $\bar{r}_2$  have identical variances and a known correlation coefficient  $\rho$ . There is a risk-free asset with rate of return  $r_f$ .

- (a) Find an expression for the optimal (Markowitz) weights for the two assets.
- (b) For the parameters  $\bar{r}_1 = .10, \bar{r}_2 = .08, r_f = .05, \rho = .6$ , find the weight of asset 1.

12. (Equivalence) Show that for equation (6.10) all rates of return  $r$  can be transformed by the linear relation  $R = ar + b$ ,  $a > 0$  and that equation (6.10) will still hold for the  $R$ 's (although the  $v_i$ 's will need to adjust as well).

13. (Coin flip) Two risky assets are derived by a single flip of a coin. For asset A, a “heads” outcome pays \$4.00, while a “tails” outcome pays \$0.00. For asset B, the corresponding



**FIGURE 6.15** Payoffs from coin flip.

payments are \$3.00 and \$1.00. The cost to invest in either asset is \$1.00. See Figure 6.15. Find the efficient portfolio made from the two assets A and B. [Hint: Use Exercise 12.]

## References

Mean-variance portfolio theory was initially devised by Markowitz [1–4]. Other important developments were presented in [5–8]. The one-fund argument is due to Tobin [9]. For comprehensive textbook presentations, see [10] and the other general investment textbooks listed as references for Chapter 2.

1. Markowitz, H. M. (1952), “Portfolio Selection,” *Journal of Finance*, **7**, no. 1, 77–91.
2. Markowitz, H. M. (1956), “The Optimization of a Quadratic Function Subject to Linear Constraints,” *Naval Research Logistics Quarterly*, **3**, nos. 1–2, 111–133.
3. Markowitz, H. M. (1987), *Portfolio Selection*, Wiley, New York.
4. Markowitz, H. M. (1987), *Mean–Variance Analysis in Portfolio Choice and Capital Markets*, Basil Blackwell, New York.
5. Hester, D. D., and J. Tobin (1967), *Risk Aversion and Portfolio Choice*, Wiley, New York.
6. Fama, E. F. (1976), *Foundations of Finance*, Basic Books, New York.
7. Sharpe, W. F. (1967), “Portfolio Analysis,” *Journal of Financial and Quantitative Analysis*, **2**, 76–84.
8. Levy, H. (1979), “Does Diversification Always Pay?” *TIMS Studies in Management Science*.
9. Tobin, J. (1958), “Liquidity Preference as Behavior Toward Risk,” *Review of Economic Studies*, **26**, February, 65–86.
10. Elton, E. J., M. J. Gruber, S. J. Brown, and W. N. Goetz, (2009), *Modern Portfolio Theory and Investment Analysis*, 8th ed., Wiley, New York.

# THE CAPITAL ASSET PRICING MODEL

**T**wo main problem types dominate the discipline of investment science. The first is to determine the best course of action in an investment situation. Problems of this type include how to devise the best portfolio, how to devise the optimal strategy for managing an investment, how to select from a group of potential investment projects, and so forth. Several examples of such problems were treated in Part 1 of this book. The second type of problem is to determine the correct, arbitrage-free, fair, or equilibrium price of an asset. We saw examples of this in Part 1 as well, such as the formula for the correct price of a bond in terms of the term structure of interest rates, and the formula for the appropriate value of a firm.

This chapter concentrates mainly on the pricing issue. It deduces the correct price of a risky asset within the framework of the mean–variance setting. The result is the **capital asset pricing model** (CAPM) developed primarily by Sharpe, Lintner, and Mossin, which follows logically from the Markowitz mean–variance portfolio theory described in the previous chapter. Later in this chapter we discuss how this result can be applied to investment decision problems.

## 7.1 Market Equilibrium

Suppose that everyone is a mean–variance optimizer as described in the previous chapter. Suppose further that everyone agrees on the probabilistic structure of assets; that is, everyone assigns to the returns of assets the same mean values, the same variances, and the same covariances. Furthermore, assume that there is a unique

risk-free rate of borrowing and lending that is available to all, and that there are no transactions costs. With these assumptions what will happen?

From the one-fund theorem we know that everyone will purchase a single fund of risky assets, and they may, in addition, borrow or lend at the risk-free rate. Furthermore, since everyone uses the same means, variances, and covariances, everyone will use the same risky fund. The mix of these two assets, the risky fund and the risk-free asset, will likely vary across individuals according to their individual tastes for risk. Some will seek to avoid risk and will, accordingly, have a high percentage of the risk-free asset in their portfolios; others, who are more aggressive, will have a high percentage of the risky fund. However, every individual will form a portfolio that is a mix of the risk-free asset and the single, risky *one fund*. Hence the *one fund* in the theorem is really the *only fund* that is used.

If everyone purchases the same fund of risky assets, what must that fund be? The answer to this question is the key insight underlying the CAPM. A bit of reflection reveals that the answer is that this fund must equal the **market portfolio**. The market portfolio is the summation of all assets. In the world of equity securities, it is the totality of shares of AAPL, DOW, PEP, and so forth. If everyone buys just one fund, and their purchases add up to the market, then that one fund must be the market as well; that is, it must contain shares of every stock in proportion to that stock's representation in the entire market.

An asset's weight in a portfolio is defined as the proportion of portfolio capital that is allocated to that asset. Hence the weight of an asset in the market portfolio is equal to the proportion of that asset's total capital value to the total market capital value. These weights are termed **capitalization weights**. It is these weights that we usually denote by  $w_i$ . In other words, the  $w_i$ 's of the market portfolio are the capitalization weights of the assets.

The exact definition of the market portfolio is illustrated as follows. Suppose there are only three stocks in the market: Jazz, Inc., Classical, Inc., and Rock, Inc. Their outstanding shares and prices are shown in Table 7.1. The market weights are proportional to the total market capitalization, not to the number of shares.

**TABLE 7.1**  
**MARKET CAPITALIZATION WEIGHTS**

Security	Shares outstanding	Relative shares in market	Price	Capitalization	Weight in market
Jazz, Inc.	10,000	1/8	\$6.00	\$60,000	3/20
Classical, Inc.	30,000	3/8	\$4.00	\$120,000	3/10
Rock, Inc.	40,000	1/2	\$5.50	\$220,000	11/20
Total	80,000	1		\$400,000	1

*The percentage of shares of a stock in the market portfolio is a share-weighted proportion of total shares. These percentages are not the market portfolio weights. The market portfolio weight of a stock is proportional to capitalization. If the price of an asset changes, the share proportions do not change, but the capitalization weights do change.*

In the situation where everyone follows the mean–variance methodology with the same estimates of parameters, we know that the efficient fund of risky assets will be the market portfolio. Hence under these assumptions there is no need for us to formulate the mean–variance problem, to estimate the underlying parameters, or to solve the system of equations that define the optimal portfolio. We know that the optimal portfolio will turn out to be the market portfolio.

How does this happen? How can it be that we solve the problem even without knowing the required data? The answer is based on an **equilibrium** argument. If everyone else (or at least a large number of people) solves the problem, we do not need to. It works like this: The return on an asset depends on both its initial price and its final price. The other investors solve the mean–variance portfolio problem using their common estimates, and they place orders in the market to acquire their portfolios. If the orders placed do not match what is available, the prices must change. The prices of assets under heavy demand will increase; the prices of assets under light demand will decrease. These price changes affect the estimates of asset returns directly, and hence investors will recalculate their optimal portfolios. This process continues until demand exactly matches supply; that is, it continues until there is equilibrium.

In the idealized world, where every investor is a mean–variance investor and all have the same estimates, everyone buys the same portfolio, and that must be equal to the market portfolio. In other words, prices adjust to drive the market to efficiency. Then after other people have made the adjustments, we can be sure that the efficient portfolio is the market portfolio, so we need not make any calculations.

This theory of equilibrium is usually applied to assets that are traded repeatedly over time, such as the stock market. In this case it is argued that individuals adjust their return estimates slowly, and only make a series of minor adjustments to their calculations rather than solving the entire portfolio optimization problem at one time.

Finally, in such equilibrium models it is argued that the appropriate equilibrium need be calculated by only a few devoted (and energetic) individuals. They move prices around to the proper value, and other investors follow their lead by purchasing the market portfolio.

These arguments about the equilibrium process all have a degree of plausibility, and all have weaknesses. Deeper analysis can be carried out, but for our purposes we will merely consider that equilibrium occurs. Hence the ultimate conclusion of the mean–variance approach is that the *one fund* must be the market portfolio.

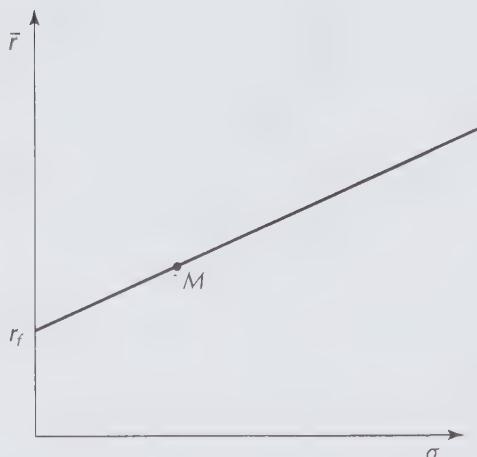
## 7.2 The Capital Market Line

Given the preceding conclusion that the single efficient fund of risky assets is the market portfolio, we can label this fund on the  $\bar{r} - \sigma$  diagram with an *M* for *market*. The efficient set therefore consists of a single straight line, emanating from the risk-free point and passing through the market portfolio. This line, shown in Figure 7.1, is called the **capital market line**.

This line shows the relation between the expected rate of return and the risk of return (as measured by the standard deviation) for efficient assets or portfolios of

**FIGURE 7.1 Capital market line.**

Efficient assets must all lie on the line determined by the risk-free rate and the market portfolio.



assets. It is also referred to as a pricing line, since prices should adjust so that efficient assets fall on this line.

The line has great intuitive appeal. It states that as risk increases, the corresponding expected rate of return must also increase. Furthermore, this relationship can be described by a straight line if risk is measured by standard deviation. In mathematical terms the capital market line states that

$$\bar{r} = r_f + \frac{\bar{r}_M - r_f}{\sigma_M} \sigma, \quad (7.1)$$

where  $\bar{r}_M$  and  $\sigma_M$  are the expected value and the standard deviation of the market rate of return and  $\bar{r}$  and  $\sigma$  are the expected value and the standard deviation of the rate of return of an arbitrary efficient asset.

The slope of the capital market line is  $K = (\bar{r}_M - r_f)/\sigma_M$ , and this value is frequently called the **price of risk**. It tells by how much the expected rate of return of a portfolio must increase if the standard deviation of that rate increases by one unit.

**Example 7.1 (The impatient investor)** Mr. Smith is young and impatient. He notes that the risk-free rate is only 6% and the market portfolio of risky assets has an expected return of 12% and a standard deviation of 15%. He figures that it would take about 60 years for his \$1,000.00 nest egg to increase to \$1 million if it earned the market rate of return. He can't wait that long. He wants that \$1 million in 10 years.

Mr. Smith easily determines that he must attain an average rate of return of about 100% per year to achieve his goal (since  $\$1,000 \times 2^{10} = \$1,048,000$ ). Correspondingly, his yearly standard deviation according to the capital market line would be the value of  $\sigma$  satisfying

$$1.0 = .06 + \frac{12. - .06}{15} \sigma,$$

or  $\sigma = 10$ . This corresponds to  $\sigma = 1,000\%$ . So this young man is certainly not guaranteed success (even if he could borrow the amount required to move far beyond the market on the capital market line).

**Example 7.2 (An oil venture)** Consider an oil drilling venture. The price of a share of this venture is \$875. It is expected to yield the equivalent of \$1,000 after 1 year, but due to high uncertainty about how much oil is at the drilling site, the standard deviation of the return is  $\sigma = 40\%$ . Currently the risk-free rate is 10%. The expected rate of return on the market portfolio is 17%, and the standard deviation of this rate is 12%.

Let us see how this venture compares with assets on the capital market line. Given the level of  $\sigma$ , the expected rate of return predicted by the capital market line is

$$\bar{r} = .10 + \frac{.17 - .10}{12} \cdot .40 = 33\%.$$

However, the actual expected rate of return is only  $\bar{r} = 1,000/875 - 1 = 14\%$ . Therefore the point representing the oil venture lies well below the capital market line. (This does *not* mean that the venture is necessarily a poor one, as we shall see later, but it certainly does not, by itself, constitute an efficient portfolio.)

## 7.3 The Pricing Model

The capital market line relates the expected rate of return of an efficient portfolio to its standard deviation, but it does not show how the expected rate of return of an individual asset relates to its individual risk. This relation is expressed by the capital asset pricing model.

We state this major result as a theorem. The reader may wish merely to glance over the proof at first reading since it is a bit involved. We shall discuss the implications of the result following the proof.

**The capital asset pricing model (CAPM)** *If the market portfolio  $M$  is efficient, the expected return  $\bar{r}_i$  of any asset  $i$  satisfies*

$$\bar{r}_i - r_f = \beta_i(\bar{r}_M - r_f), \quad (7.2)$$

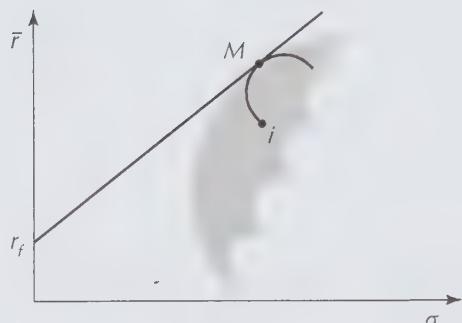
where

$$\beta_i = \frac{\sigma_{iM}}{\sigma_M^2}. \quad (7.3)$$

**Proof:** For any  $\alpha$  consider the portfolio consisting of a portion  $\alpha$  invested in asset  $i$  and a portion  $1 - \alpha$  invested in the market portfolio  $M$ . (We allow  $\alpha < 0$ , which corresponds to borrowing at the risk-free rate.) The expected rate of return of this portfolio is

$$\bar{r}_\alpha = \alpha \bar{r}_i + (1 - \alpha) \bar{r}_M,$$

**FIGURE 7.2 Portfolio curve.** The family of portfolios traces out a curve on the diagram. This curve cannot cross the capital market line, and hence must be tangent to that line.



and the standard deviation of the rate of return is

$$\sigma_\alpha = [\alpha^2 \sigma_i^2 + 2\alpha(1-\alpha)\sigma_{iM} + (1-\alpha)^2 \sigma_M^2]^{1/2}.$$

As  $\alpha$  varies, these values trace out a curve in the  $\bar{r} - \sigma$  diagram, as shown in Figure 7.2. In particular,  $\alpha = 0$  corresponds to the market portfolio  $M$ . This curve cannot cross the capital market line. If it did, the portfolio corresponding to a point above the capital market line would violate the very definition of the capital market line as being the efficient boundary of the feasible set. Hence as  $\alpha$  passes through zero, the curve must be tangent to the capital market line at  $M$ . This tangency is the condition that we exploit to derive the formula.

The tangency condition can be translated into the condition that the slope of the curve is equal to the slope of the capital market line at the point  $M$ . To set up this condition we need to calculate a few derivatives.

First we have

$$\begin{aligned}\frac{d\bar{r}_\alpha}{d\alpha} &= \bar{r}_i - \bar{r}_M \\ \frac{d\sigma_\alpha}{d\alpha} &= \frac{\alpha\sigma_i^2 + (1-2\alpha)\sigma_{iM} + (\alpha-1)\sigma_M^2}{\sigma_\alpha}.\end{aligned}$$

Thus,

$$\left. \frac{d\sigma_\alpha}{d\alpha} \right|_{\alpha=0} = \frac{\sigma_{iM} - \sigma_M^2}{\sigma_M}.$$

We then use the relation

$$\frac{d\bar{r}_\alpha}{d\sigma_\alpha} = \frac{d\bar{r}_\alpha/d\alpha}{d\sigma_\alpha/d\alpha}$$

to obtain

$$\left. \frac{d\bar{r}_\alpha}{d\sigma_\alpha} \right|_{\alpha=0} = \frac{(\bar{r}_i - \bar{r}_M)\sigma_M}{\sigma_{iM} - \sigma_M^2}.$$

This slope must equal the slope of the capital market line. Hence,

$$\frac{(\bar{r}_i - \bar{r}_M)\sigma_M}{\sigma_{iM} - \sigma_M^2} = \frac{\bar{r}_M - r_f}{\sigma_M}.$$

We now just solve for  $\bar{r}_i$ , obtaining the final result

$$\bar{r}_i = r_f + \left( \frac{\bar{r}_M - r_f}{\sigma_M^2} \right) \sigma_{iM} = r_f + \beta_i (\bar{r}_M - r_f).$$

This is clearly equivalent to the stated formula. ■

The value  $\beta_i$  is referred to as the **beta** of an asset. When the asset is fixed in a discussion, we often just write beta without a subscript— $\beta$ . An asset's beta is all that need be known about the asset's risk characteristics to use the CAPM formula.

The value  $\bar{r}_i - r_f$  is termed the **expected excess rate of return** of asset  $i$ ; it is the amount by which the rate of return is expected to exceed the risk-free rate. Likewise,  $\bar{r}_M - r_f$  is the expected excess rate of return of the market portfolio. In terms of these expected excess rates of return, the CAPM says that the expected excess rate of return of an asset is proportional to the expected excess rate of return of the market portfolio, and the proportionality factor is  $\beta$ . So with  $r_f$  taken as a base point, the expected returns of a particular asset and of the market above that base are proportional.

An alternative interpretation of the CAPM formula is based on the fact that  $\beta$  is a normalized version of the covariance of the asset with the market portfolio. Hence the CAPM formula states that the expected excess rate of return of an asset is directly proportional to its covariance with the market. It is this covariance that determines the expected excess rate of return.

To gain insight into this result, let us consider some extreme cases. Suppose, first, that the asset is completely *uncorrelated* with the market; that is,  $\beta = 0$ . Then, according to the CAPM, we have  $\bar{r} = r_f$ . This is perhaps at first sight a surprising result. It states that even if the asset is very risky (with large  $\sigma$ ), the expected rate of return will be that of the risk-free asset—there is no premium for risk. The reason for this is that the risk associated with an asset that is uncorrelated with the market can be diversified away. If we had many such assets, each uncorrelated with the others and with the market, we could purchase small amounts of each of them, and the resulting total variance would be small. Since the final composite return would have small variance, the corresponding expected rate of return should be close to  $r_f$ .

Even more extreme is an asset with a negative value of  $\beta$ . In that case  $\bar{r} < r_f$ ; that is, even though the asset may have very high risk (as measured by its  $\sigma$ ), its expected rate of return should be even less than the risk-free rate. The reason is that such an asset reduces the overall portfolio risk when it is combined with the market. Investors are therefore willing to accept the lower expected value for this risk-reducing potential. Such assets provide a form of insurance. They do well when everything else does poorly.

The CAPM changes our concept of the risk of an asset from that of  $\sigma$  to that of  $\beta$ . It is still true that, overall, we measure the risk of a portfolio in terms of  $\sigma$ , but this does not translate into a concern for the  $\sigma$ 's of individual assets. For those, the proper measure is their  $\beta$ 's.

**Example 7.3 (A simple calculation)** We illustrate how simple it is to use the CAPM formula to calculate an expected rate of return. Let the risk-free rate be  $r_f = 8\%$ .

Suppose the rate of return of the market has an expected value of 12% and a standard deviation of 15%.

Now consider an asset that has covariance of .045 with the market. Then we find  $\beta = .045/(.15)^2 = 2.0$ . The expected return of the asset is  $\bar{r} = .08 + 2 \times (.12 - .08) = .16 = 16\%$ .

## Betas of Common Stocks

The concept of beta is well established in the financial community, and it is referred to frequently in technical discussions about particular stocks. Beta values are estimated by various financial service organizations. Typically, these estimates are formed by using a record of past stock values (usually about 6 or 18 months of weekly values) and computing, from the data, average values of returns, products of returns, and squares of returns in order to approximate expected returns, covariances, and variances. The beta values so obtained drift around somewhat over time, but unless there are drastic changes in a company's situation, its beta tends to be relatively stable.

Table 7.2 lists some well-known U.S. companies and their corresponding beta ( $\beta$ ) and volatility ( $\sigma$ ) values as estimated at a particular date. Try scanning the list and see if the values given support your intuitive impression of the company's market properties. Generally speaking, we expect aggressive companies or highly leveraged companies to have high betas, whereas conservative companies whose performance is unrelated to the general market behavior are expected to have low betas. Also, we expect that companies in the same business will have similar, but not identical, beta values. Compare, for instance, Coca-Cola with Pepsi, or Intel with Texas Instruments.

## Beta of a Portfolio

It is easy to calculate the overall beta of a portfolio in terms of the betas of the individual assets in the portfolio. Suppose, for example, that a portfolio contains  $n$  assets with the weights  $w_1, w_2, \dots, w_n$ . The rate of return of the portfolio is  $r = \sum_{i=1}^n w_i r_i$ . Hence  $\text{cov}(r, r_M) = \sum_{i=1}^n w_i \text{cov}(r_i, r_M)$ . It follows immediately that

$$\beta_p = \sum_{i=1}^n w_i \beta_i. \quad (7.4)$$

In other words, the portfolio beta is just the weighted average of the betas of the individual assets in the portfolio, with the weights being identical to those that define the portfolio.

## 7.4 The Security Market Line

The CAPM formula can be expressed in graphical form by regarding the formula as a linear relationship. This relationship is termed the **security market line**. Two versions are shown in Figure 7.3.

**TABLE 7.2**  
**SOME U.S. COMPANIES: THEIR BETAS AND SIGMAS**

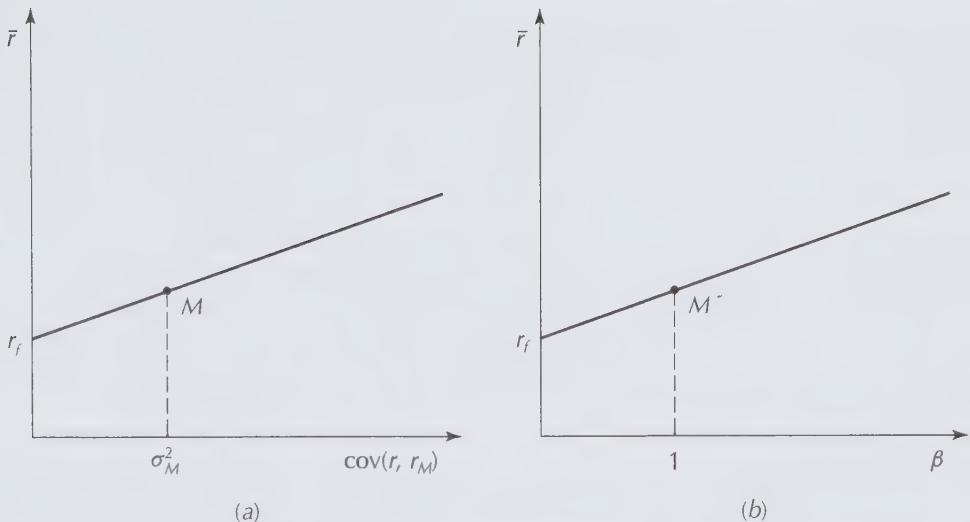
Ticker sym	Company name	Beta	Sigma(%)
AA	Alcoa	1.85	55.14
AAPL	Apple	1.22	39.34
ALL	Allstate	1.28	45.65
AMZN	Amazon	1.17	49.20
AVP	Avon Products	1.01	36.24
BA	Boeing	1.21	33.93
CAT	Caterpillar	1.61	39.12
COST	COSCO Wholesale	0.67	27.46
CVX	Chevron	0.91	32.44
DELL	Dell	0.97	40.18
DOW	Dow Chemical	1.30	44.53
F	Ford Motor	2.04	60.47
FDX	Fedex	1.10	36.21
GE	General Electric	1.40	39.08
GOOG	Google	0.96	34.44
GS	Goldman Sachs Gro	1.37	49.83
INTC	Intel	1.17	35.08
KFT	Kraft Foods	0.49	23.00
KO	Coca-Cola	0.45	21.50
MCD	McDonalds	0.51	22.67
MS	Morgan Stanley	1.88	71.43
MSFT	Microsoft	0.86	31.96
PEP	Pepsico	0.44	20.62
T	AT&T	0.69	27.62
TWX	Time Warner	1.12	36.91
TXN	Texas Instruments	1.00	33.13
UNP	Union Pacific	1.29	36.41
WAG	Walgreen	0.64	29.93
WMT	Wal-Mart Stores	0.49	22.34
XOM	Exxon Mobil	0.69	30.29

Source: ABG Analytics, April 29, 2011.

Both graphs show the linear variation of  $\bar{r}$ . The first expresses it in covariance form, with  $\text{cov}(r, r_M)$  being the horizontal axis. The market portfolio corresponds to the point  $\sigma_M^2$  on this axis. The second graph shows the relation in beta form, with beta being the horizontal axis. In this case the market corresponds to the point  $\beta = 1$ .

Both of these lines highlight the essence of the CAPM formula. Under the equilibrium conditions assumed by the CAPM, any asset should fall on the security market line.

The security market line expresses the risk–reward structure of assets according to the CAPM, and emphasizes that the risk of an asset is a function of its covariance with the market or, equivalently, a function of its beta.



**FIGURE 7.3 Security market line.** The expected rate of return increases linearly as the covariance with the market increases or, equivalently, as  $\beta$  increases.

## Systematic Risk

The CAPM implies a special structural property for the return of an asset, and this property provides further insight as to why beta is the most important measure of risk. To develop this result we write the (random) rate of return of asset  $i$  as

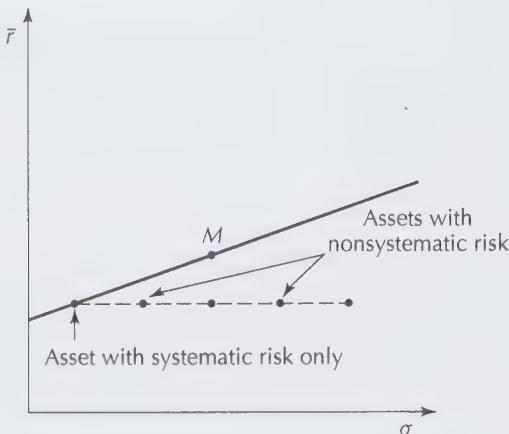
$$r_i = r_f + \beta_i(r_M - r_f) + \varepsilon_i. \quad (7.5)$$

This is just an arbitrary equation at this point. The random variable  $\varepsilon_i$  is chosen to make it true. However, the CAPM formula tells us several things about  $\varepsilon_i$ .

First, taking the expected value of (7.5), the CAPM says that  $E(\varepsilon_i) = 0$ . Second, taking the correlation of (7.5) with  $r_M$  (and using the definition of  $\beta_i$ ), we find  $\text{cov}(\varepsilon_i, \sigma_M) = 0$ . We can therefore write

$$\sigma_i^2 = \beta_i^2 \sigma_M^2 + \text{var}(\varepsilon_i)$$

and we see that  $\sigma_i^2$  is the sum of two parts. The first part,  $\beta_i^2 \sigma_M^2$ , is termed the **systematic risk**. This is the risk associated with the market as a whole. This risk cannot be reduced by diversification because every asset with nonzero beta contains this risk. The second part,  $\text{var}(\varepsilon_i)$ , is termed the **nonsystematic, idiosyncratic, or specific risk**. This risk is uncorrelated with the market and can be reduced by diversification. It is the systematic (or nondiversifiable) risk, measured by beta, that is most important, since it directly combines with the systematic risk of other assets.



**FIGURE 7.4 Systematic and nonsystematic risk.** An asset on the capital market line has only systematic risk. Assets with nonsystematic risk fall to the right of the capital market line.

Consider an asset on the capital market line<sup>1</sup> with a value of  $\beta$ . The standard deviation of this asset is  $\beta\sigma_M$ . It has only systematic risk; there is no nonsystematic risk. This asset has an expected rate of return equal to  $\bar{r} = r_f + \beta(\bar{r}_M - r_f)$ . Now consider a whole group of other assets, all with the same value of  $\beta$ . According to CAPM, these all have the same expected rate of return, equal to  $\bar{r}$ . However, if these assets carry nonsystematic risk, they will not fall on the capital market line. Indeed, as the nonsystematic risk increases, the points on the  $\bar{r} - \sigma$  plane representing these assets drift to the right, as shown in Figure 7.4. The horizontal distance of a point from the capital market line is therefore a measure of the nonsystematic risk.

## 7.5 Investment Implications

The question of interest for the investor is: Can the CAPM help with investment decisions? There is no simple answer to this question.

The CAPM states (or assumes), based on an equilibrium argument, that the solution to the Markowitz problem is that the market portfolio is the *one fund* (and *only fund*) of risky assets that anyone need hold. This fund is supplemented only by the risk-free asset. The investment recommendation that follows this argument is that an investor should simply purchase the market portfolio. That is, ideally, an investor should purchase a little bit of every asset that is available, with the proportions determined by the relative amounts that are issued in the market as a whole. If the world of equity securities is taken as the set of available assets, then each person should purchase some shares in every available stock, in proportion to the stocks' monetary share of the total of all stocks outstanding. It is not necessary to go to

<sup>1</sup> Of course, to be exactly on the line, the asset must be equivalent to a combination of the market portfolio and the risk-free asset.

the trouble of analyzing individual issues and computing a Markowitz solution. Just apportion total investment between the market portfolio and the risk-free asset.

Since it would be rather cumbersome for an individual to assemble the market portfolio, mutual funds have been designed to match the market portfolio closely. These funds are termed **index funds**, since they usually attempt to duplicate the portfolio of a major stock market index, such as the *Standard & Poor's 500* (S&P 500), an average of 500 stocks that as a group is thought to be representative of the market as a whole. Other indices use even larger numbers of stocks. A CAPM purist (that is, one who fully accepts the CAPM theory as applied to publicly traded securities) could just purchase one of these index funds (to serve as the *one fund*) as well as some risk-free securities such as U.S. Treasury bills.

Some people believe that they can do better than blindly purchasing the market portfolio. The CAPM, after all, assumes that everyone has identical information about the (uncertain) returns of all assets. Clearly, this is not the case. If someone believes that he or she possesses superior information, then presumably that person could form a portfolio that would outperform the market. We return to this issue in Chapter 9, where questions concerning data and information are explicitly addressed. It is shown there that it is not at all easy to obtain accurate data for use in a Markowitz model, and hence the solution computed from such a model is likely to be somewhat nonsensical. For now we just state that the best designs seem to be those formulated as deviations or extensions of the basic CAPM idea, rather than as bold new beginnings. In other words, in constructing a portfolio, one probably should begin with the market portfolio and alter it systematically, rather than attempting to solve the full Markowitz problem from scratch.

One area where the CAPM approach has direct application is in the analysis of assets that do not have well-established market prices. In this case the CAPM can be used to find a *reasonable* price. An important class of problems of this type are the project evaluation problems (variations of capital budgeting problems) that arise in firms. This application is considered explicitly in Sections 7.8–7.10.

## 7.6 Performance Evaluation

The CAPM theory can be used to evaluate the performance of an investment portfolio, and indeed it is now common practice to evaluate many institutional portfolios (such as pension funds and mutual funds) using the CAPM framework. We shall present the main ideas by going through a simple hypothetical example. The primary purpose of this section, however, is to use these performance measure ideas to illustrate the CAPM.

**Example 7.4 (ABC fund analysis)** The ABC mutual fund has the 10-year record of rates of return shown in the column labeled ABC in Table 7.3. We would like to evaluate this fund's performance in terms of mean-variance portfolio theory and the CAPM. Is it a good fund that we could recommend? Can it serve as the *one fund* for a prudent mean-variance investor?

**TABLE 7.3**  
**ABC FUND PERFORMANCE**

<b>Year</b>	<b>Rate of return percentages</b>		
	<b>ABC</b>	<b>S&amp;P</b>	<b>T-bills</b>
1	14	12	7
2	10	7	7.5
3	19	20	7.7
4	-8	-2	7.5
5	23	12	8.5
6	28	23	8
7	20	17	7.3
8	14	20	7
9	-9	-5	7.5
10	19	16	8
Average	13	12	7.6
Standard deviation	12.4	9.4	.5
Geometric mean	12.3	11.6	7.6
Cov(ABC, S&P)	.0107		
Beta	1.20375	1	
Jensen	0.00104	0.00000	
Sharpe	0.43577	0.46669	

The top part of the table shows the rate of return achieved by ABC, S&P 500, and T-bills over a 10-year period. The lower portion shows the Jensen and Sharpe indices.

**Step 1.** We begin our analysis by computing the three quantities shown in Table 7.3 below the given return data: the average rate of return, the standard deviation of the rate as implied by the 10 samples, and the geometric mean rate of return. These quantities are estimates based on the available data.

In general, given  $r_i, i = 1, 2, \dots, n$ , the average rate of return is

$$\hat{r} = \frac{1}{n} \sum_{i=1}^n r_i$$

and this serves as an estimate of the true expected return  $\bar{r}$ . The average variance is<sup>2</sup>

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (r_i - \hat{r})^2$$

and the estimate  $s$  of the standard deviation is the square root of that. It is also useful to calculate the geometric mean rate of return, which is

$$\mu = [(1 + r_1)(1 + r_2) \cdots (1 + r_n)]^{1/n} - 1.$$

<sup>2</sup> The reason that  $n - 1$  is used in the denominator instead of  $n$  is discussed in Chapter 9.

This measures the actual rate of return over the  $n$  years, accounting for compounding. This value will generally be somewhat lower than the average rate of return.

**Step 2.** Next we obtain data on both the market portfolio and the risk-free rate of return over the 10-year period. We use the *Standard & Poor's 500* stock average and the 1-year Treasury bill rate, respectively. These are shown in Table 7.3. We calculate average rates of return and standard deviations of these by the same method as for ABC. We also calculate an estimate of the covariance of the ABC fund with the S&P 500 by using the estimate

$$\text{cov}(r, r_M) = \frac{1}{n-1} \sum_{i=1}^n (r_i - \hat{r})(r_{Mi} - \hat{r}_M).$$

We then calculate beta from the standard formula,

$$\beta = \frac{\text{cov}(r, r_M)}{\text{var}(r_M)}.$$

This gives us enough information to carry out an interesting analysis.

**Step 3.** (The Jensen index) We write the formula

$$\hat{r} - r_f = J + \beta(\hat{r}_M - r_f).$$

This looks like the CAPM pricing formula (7.2), except that we have replaced expected rates of return by measured average returns (for that is the best that can be done in this situation), and we have added an error term  $J$ . The  $J$  here stands for **Jensen's index**.

According to the CAPM, the value of  $J$  should be zero when true expected returns are used. Hence  $J$  measures, approximately, how much the performance of ABC has deviated from the theoretical value of zero. A positive value of  $J$  presumably implies that the fund did better than the CAPM prediction (but of course we recognize that approximations are introduced by the use of a finite amount of data to estimate the important quantities).

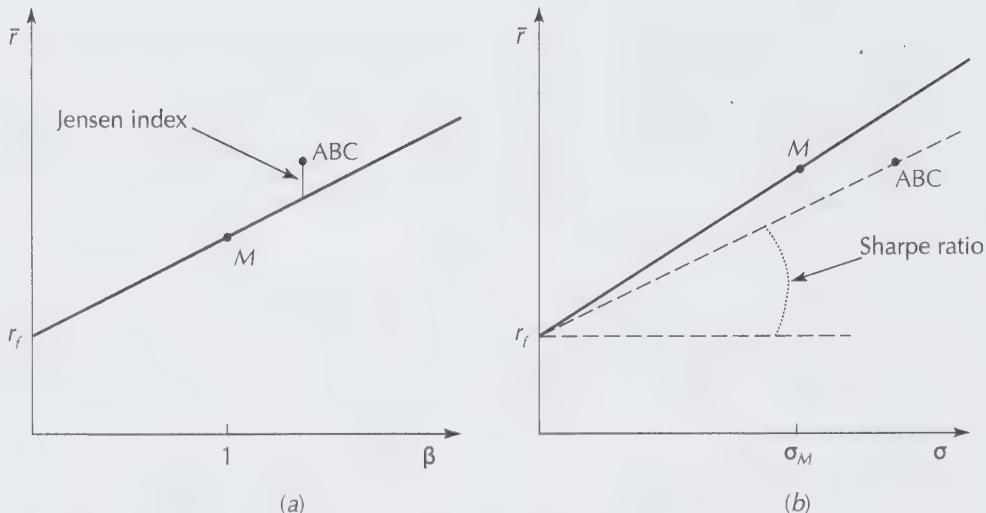
The Jensen index can be indicated on the security market line, as shown in Figure 7.5(a). For ABC, we find that indeed  $J > 0$ , and hence we might conclude that ABC is an excellent fund. But is this really a correct inference?

Aside from the difficulties inherent in using short histories of data this way, the inference that ABC is a good mutual fund is not entirely warranted. It is *not* clear that it can serve as the one fund of risky assets in an efficient portfolio. The fact that  $J > 0$  is nice, and may tell us that ABC is a good *asset*, but it does not say that the ABC fund is, by itself, *efficient*.<sup>3</sup>

**Step 4.** (The Sharpe ratio) In order to measure the efficiency of ABC we must see *where* it falls relative to the capital market line. Only portfolios on that line are

---

<sup>3</sup> It can be argued that the Jensen index tells us nothing about the fund, but instead is a measure of the validity of the CAPM. If the CAPM is valid, then every security (or fund) must satisfy the CAPM formula exactly, since the formula is an identity if the market portfolio is efficient. If we find a security with a nonzero Jensen index, then that is a sign that the market is not efficient. The CAPM formula is often applied to (new) financial instruments or projects that are not traded and hence not part of the market portfolio. In this case, the Jensen index can be a useful measure.



**FIGURE 7.5 Performance indices for ABC.** The Jensen index measures the height above the security market line; the Sharpe ratio measures the angle in the  $\bar{r} - \sigma$  plane.

efficient. We do this by writing the formula

$$\hat{\bar{r}} - r_f = S\sigma.$$

The value of  $S$  is the slope of the line drawn between the risk-free point and the ABC point on the  $\bar{r} - \sigma$  diagram. The  $S$  stands for **Sharpe ratio**. For ABC we find  $S = .43577$ . This must be compared with the corresponding value for the market—represented by the S&P 500. We find the value for the S&P 500 is  $S = .46669$ . The situation is shown in Figure 7.5(b). Clearly ABC is not efficient, at least as revealed by the available data.

We conclude that ABC may be worth holding in a portfolio. By itself it is not quite efficient, so it would be necessary to supplement this fund with other assets or funds to achieve efficiency. Or, to attain efficiency, an investor could simply invest in a broad-based fund instead of the ABC fund.

## 7.7 CAPM as a Pricing Formula

The CAPM is a **pricing model**. However, the standard CAPM formula does not contain prices explicitly—only expected rates of return. To see why the CAPM is called a pricing model we must go back to the definition of return.

Suppose that an asset is purchased at price  $P$  and later sold at price  $Q$ . The rate of return is then  $r = (Q - P)/P$ . Here  $P$  is known and  $Q$  is random. Putting this in the

CAPM formula, we have

$$\frac{\bar{Q} - P}{P} = r_f + \beta(\bar{r}_M - r_f).$$

Solving for  $P$  we obtain

$$P = \frac{\bar{Q}}{1 + r_f + \beta(\bar{r}_M - r_f)}.$$

This gives the price of the asset according to the CAPM. We highlight this important result:

**Pricing form of the CAPM** *The price  $P$  of an asset with payoff  $\bar{Q}$  is*

$$P = \frac{\bar{Q}}{1 + r_f + \beta(\bar{r}_M - r_f)}, \quad (7.6)$$

*where  $\beta$  is the beta of the asset.*

This pricing formula has a form that very nicely generalizes the familiar discounting formula for deterministic situations. In the deterministic case, it is appropriate to discount the future payment at the interest rate  $r_f$ , using a factor of  $1/(1+r_f)$ . In the random case the appropriate interest rate is  $r_f + \beta(\bar{r}_M - r_f)$ , which can be regarded as a risk-adjusted interest rate.

**Example 7.5 (The price is right)** Gavin Jones is good at math, but his friends tell him that he doesn't always see the *big picture*. Right now, Gavin is thinking about investing in a mutual fund. This fund invests 10% of its funds at the risk-free rate of 7% and the remaining 90% in a widely diversified portfolio that closely approximates the market portfolio, which has an expected rate of return equal to 15%. One share of the mutual fund represents \$100 of assets in the fund. Having just studied the CAPM, Gavin wants to know how much such a share should cost.

Gavin figures out that the beta of the fund must be .90. The value of a share after 1 year is expected to be  $10 \times 1.07 + 90 \times 1.15 = 114.20$ . Hence, according to (7.6),

$$P = \frac{114.20}{1.07 + .90 \times (.15 - .07)} = \$100.$$

Yes, the price of a share will be equal to the value of the funds it represents. Gavin is reassured (but suspects he could have concluded that more simply).

**Example 7.6 (The oil venture)** Consider again, as in Example 7.2, the possibility of investing in a share of a certain oil well that will produce a payoff that is random because of the uncertainty associated with whether or not there is oil at that site and because of the uncertainty in future oil prices. The expected payoff is \$1,000 and the standard deviation of return is a relatively high 40%. The beta of the asset is  $\beta = .6$ , which is relatively low because, although the uncertainty in return due to oil prices is correlated with the market portfolio, the uncertainty associated with exploration is not. The risk-free rate is  $r_f = 10\%$ , and the expected return on the market portfolio is

.17. What is the value of this share of the oil venture, based on CAPM? (Recall that earlier it was stated that the offered price was \$875.) We have immediately

$$P = \frac{\$1,000}{1.10 + .6(.17 - .10)} = \$876,$$

and  $\sigma$  does not enter the calculation.

The venture may be quite risky in the traditional sense of having a high standard deviation associated with its return. But, nevertheless, it is fairly priced because of the relatively low beta.

## Linearity of Pricing and the Certainty Equivalent Form

We now discuss a very important property of the pricing formula—namely, that it is **linear**. This means that the price of the sum of two assets is the sum of their prices, and the price of a multiple of an asset is the same multiple of the price. This is really quite startling because the formula does not *look* linear at all (at least for sums). For example, if

$$P_1 = \frac{\bar{Q}_1}{1 + r_f + \beta_1(\bar{r}_M - r_f)}, \quad P_2 = \frac{\bar{Q}_2}{1 + r_f + \beta_2(\bar{r}_M - r_f)},$$

it does not seem obvious that

$$P_1 + P_2 = \frac{\bar{Q}_1 + \bar{Q}_2}{1 + r_f + \beta_{1+2}(\bar{r}_M - r_f)},$$

where  $\beta_{1+2}$  is the beta of a new asset, which is the sum of assets 1 and 2. Furthermore, based on our recognition that the covariance between assets is important in assessing how to use them in a portfolio, it may seem *unreasonable* that the pricing formula should be linear. We can easily take care of the first doubt by converting the formula into another form, which appears linear; then we will discuss the intuition behind the result.

The form of the CAPM pricing formula that clearly displays linearity is called the **certainty equivalent form**. Suppose that we have an asset with price  $P$  and final value  $Q$ . Here again  $P$  is known and  $Q$  is uncertain. Using the fact that  $r = Q/P - 1$ , the value of beta is

$$\beta = \frac{\text{cov}[(Q/P - 1), r_M]}{\sigma_M^2}.$$

Since constant additive terms do not affect covariance, this becomes

$$\beta = \frac{\text{cov}(Q, r_M)}{P\sigma_M^2}.$$

Substituting this into the pricing formula (7.6) and dividing by  $P$  yields

$$1 = \frac{\bar{Q}}{P(1 + r_f) + \text{cov}(Q, r_M)(\bar{r}_M - r_f)/\sigma_M^2}.$$

Finally, solving for  $P$  we obtain the following formula:

**Certainty equivalent pricing formula** *The price  $P$  of an asset with payoff  $Q$  is*

$$P = \frac{1}{1+r_f} \left[ \bar{Q} - \frac{\text{cov}(Q, r_M)(\bar{r}_M - r_f)}{\sigma_M^2} \right]. \quad (7.7)$$

The term in brackets is called the **certainty equivalent** of  $Q$ . This value is treated as a certain amount, and then the normal discount factor  $1/(1+r_f)$  is applied to obtain  $P$ . The certainty equivalent form shows clearly that the pricing formula is linear because both terms in the brackets depend linearly on  $Q$ .

The reason for linearity can be traced back to the principle of no arbitrage: if the price of the sum of two assets were not equal to the sum of the individual prices, it would be possible to make arbitrage profits. For example, if the combination asset were priced lower than the sum of the individual prices, we could buy the combination (at the low price) and sell the individual pieces (at the higher price), thereby making a profit. By doing this in large quantities, we could make arbitrarily large profits. If the reverse situation held—if the combination asset were priced higher than the sum of the two assets—we would buy the assets individually and sell the combination, again making arbitrage profits. Such arbitrage opportunities are ruled out if and only if the pricing of assets is linear. This linearity of pricing is therefore a fundamental tenet of financial theory (in the context of perfect markets), and we shall return to it frequently throughout the text.

**Example 7.7 (Gavin tries again)** Gavin Jones decides to use the certainty equivalent form of the pricing equation to calculate the share price of the mutual fund considered in Example 7.5. In this case he notes that  $\text{cov}(Q, r_M) = 90\sigma_M^2$ , where  $Q$  is the value of the fund after 1 year. Hence,

$$P = \frac{114.20 - 90 \times .08}{1.07} = \$100.$$

All is well again, according to his math.

It is frequently convenient to use payoff values rather than rates of return in the certainty equivalent form, since the  $P$ 's and  $Q$ 's are values, not rates. Thus we use  $M$  for the final market value,  $P_M$  for its price, and  $R$  for the total return of the risk-free asset. In these terms the certainty equivalent price of an asset is

$$P = \frac{1}{R} \left[ \bar{Q} - \frac{\text{cov}(Q, M)(\bar{M} - P_M R)}{\sigma_M^2} \right], \quad (7.8)$$

where now it is understood that  $\sigma_M^2$  is the the variance of  $M$  (not of  $r_M$ ).

## 7.8 Project Choice\*

A firm can use the CAPM as a basis for deciding which projects it should carry out. Suppose, for example, that a potential project requires an initial outlay of  $P$  and will generate a net amount  $Q$  after 1 year. As usual,  $P$  is known and  $Q$  is random, with expected value  $\bar{Q}$ . It is natural to define the net present value (NPV) of this project by the formula

$$\text{NPV} = -P + \frac{1}{1+r_f} \left[ \bar{Q} - \frac{\text{cov}(Q, r_M)(\bar{r}_M - r_f)}{\sigma_M^2} \right] \quad (7.9)$$

This formula is based on the certainty equivalent form of the CAPM: the first (negative) term is the initial outlay and the second term is the certainty equivalent of the final payoff.

The firm may have many different projects from which it will select a few. What criterion should the firm employ in making its selection? Extending our knowledge of the deterministic case, it seems appropriate for the firm to select the group of projects that maximize NPV. Indeed this is the advice that is normally given to firms.

How would potential investors view the situation? For them a particular firm is only one of a whole group of firms in which they may choose to invest. Investors are concerned with the overall performance of their portfolios, and only incidentally with the internal decisions of a particular firm. If investors base their investment decisions on a mean-variance criterion, they want an individual firm to operate so as to push the efficient frontier, of the entire universe of assets, as far upward and leftward as possible. This would improve the efficient frontier and hence the performance of a mean-variance efficient portfolio. Therefore potential investors will urge the management teams of firms to select projects that will shift the efficient frontier outward as far as possible, then they will invest in the efficient portfolio.

The two criteria—net present value and maximum expansion of the efficient frontier—may, it seems, be in conflict. The NPV criterion focuses on the firm itself; the efficient frontier criterion focuses on the joint effect of all firms. But really, there is no conflict. The two criteria are essentially equivalent, as stated by the following version of the harmony theorem:

**Harmony Theorem** *If a firm has a project with positive NPV, then the efficient frontier can be expanded.*

**Proof:** The entrepreneur is concerned with present value, while investors understand that this venture is but one asset in a whole portfolio. Investors are concerned with the efficient frontier.

As investments are made, the overall market adjusts so that the market portfolio  $M$  includes the new ventures.

Suppose  $P$  is the total cost of a venture (or project) and  $Q$  is the random payoff. In equilibrium the net present value based on the two-period cash flow

$(-P, Q)$  is zero. That is,

$$0 = -P + \frac{\bar{Q} - \text{cov}(Q, r_M)(r_M - r_f)/\sigma_M^2}{1 + r_f}.$$

As we know, this implies

$$\bar{r} - r_f = \beta(\bar{r}_M - r_f),$$

where  $r$  is the rate of return to the investors and  $\beta = \text{cov}(r_M, r)/\sigma_M^2$ .

Now suppose that there is a different firm project available, not in the current market, with initial cost  $P'$  and reward  $Q'$ , for which

$$-P' + \frac{\bar{Q}' - \text{cov}(Q', r_M)(r_M - r_f)/\sigma_M^2}{1 + r_f} > 0.$$

This inequality implies that

$$\bar{r}' - r_f - \text{cov}(r', r_M)(\bar{r}_M - r_f)/\sigma^2 > 0.$$

Consider portfolios with rates of return  $r_\alpha = r_M + \alpha r' - \alpha r$  for small values of  $\alpha$ , obtained by modifying the market portfolio by transferring a small increment  $\alpha$  of portfolio weight from the market to the new project. We want to show that there are such portfolios that lie above the current efficient frontier.

To show this we evaluate

$$\tan \theta \equiv \frac{\bar{r}_\alpha - r_f}{\sigma_\alpha}$$

for small  $\alpha > 0$ . Differentiation gives

$$\frac{d \tan \theta_\alpha}{d \alpha} = \frac{1}{\sigma_\alpha} \frac{d \bar{r}_\alpha}{d \alpha} - \frac{\bar{r}_\alpha - r_f}{\sigma_\alpha^2} \frac{d \sigma_\alpha}{d \alpha}.$$

Using

$$\begin{aligned} \frac{d \bar{r}_\alpha}{d \alpha} \Big|_{\alpha=0} &= \bar{r}' - \bar{r} \\ \frac{d \sigma_\alpha}{d \alpha} \Big|_{\alpha=0} &= \frac{\sigma_{MQ'} - \sigma_{MQ}}{\sigma_M}, \end{aligned}$$

we find

$$\begin{aligned} \frac{d \tan \theta_\alpha}{d \alpha} \Big|_{\alpha=0} &= \frac{\bar{r}' - \bar{r}}{\sigma_M} - \frac{\bar{r}_M - r_f}{\sigma_M^2} \frac{\sigma_{MQ'} - \sigma_{MQ}}{\sigma_M} \\ &= \frac{1}{\sigma_M} [\bar{r}' - \beta'(\bar{r}_M - r_f)] - \frac{1}{\sigma_M} [\bar{r} - \beta(\bar{r}_M - r_f)] > 0. \end{aligned}$$

The final inequality follows because the first bracketed term is greater than  $r_f$  and the second is equal to  $r_f$ . For  $\alpha$  small this means that  $\tan \theta_\alpha > \tan \theta_0$ . Hence, the efficient frontier is larger than it was originally. ■

This result shows that if the entrepreneur maximizes NPV and investors desire efficiency, the combination is in harmony.

## 7.9 Projection Pricing

Speculation about what price a new asset with random payoff should command when entered into the overall market is a common financial exercise. Potentially, there are many ways a suitable price can be assigned; accordingly, a whole range of possible prices might be inferred. A common method for pricing uses the CAPM. In this method the random payoff variable is inserted into the CAPM formula, and the resulting computed value is then considered an appropriate price. Of course the CAPM formula is “correct” only when the market portfolio is equal to the Markowitz efficient portfolio and only when the formula is restricted to assets that are already in the market or are combinations of assets in the market. It is these assets that define the Markowitz portfolio. However, the CAPM formula *will* assign a price to the asset. But what is the nature of this price?

This section and the next will investigate this situation and present methods equivalent to the CAPM but that in some situations have advantages. Furthermore, we shall develop a form of the CAPM that is quite intuitive and parallels the method of price comparison that one typically uses in everyday life.

In our development we follow the understanding that CAPM refers to the pricing formula that uses the Markowitz efficient portfolio. In general we do not assert that this is equal to the market portfolio; but it is often a useful approximation to assume that they are equivalent.

To carry out our analysis of these new procedures, it is useful to take a somewhat abstract view.

Suppose there are  $n$  marketed securities determined by their prices  $P_1, P_2, \dots, P_n$  and their payoffs  $y_1, y_2, \dots, y_n$ , which are random. Consider the  $n$  payoffs as generating a **market space**  $\mathcal{M}$  consisting of arbitrary linear combinations of the  $n$  basic payoff variables. For example, if there are two basic assets with payoffs  $y_1$  and  $y_2$ , the space  $\mathcal{M}$  is two-dimensional. The payoff  $y = 6y_1 + 3.4y_2$  is a member of this market space. Prices of elements in this space are defined by the corresponding linear combination of basic prices. For example, the price of  $y = 6y_1 + 3.4y_2$  is  $P = 6P_1 + 3.4P_2$ , according to linear pricing.<sup>4</sup>

In the CAPM framework, one of the marketed securities is risk free with payoff 1 and price  $1/R$ . The expected values of all  $y_i$ 's are known, as well as all covariances (and variances) between the  $y_i$ 's.

To find the price of a particular payoff  $y \in \mathcal{M}$  with linear pricing, we first express  $y$  in the form  $y = c_1y_1 + c_2y_2 + \dots + c_ny_n$ , where the  $y_i$ 's are given assets in

---

<sup>4</sup> For all this to work out without difficulty, it is assumed that the  $y_i$ 's are **linearly independent**, in the sense that there is no linear combination of the  $y_i$ 's that gives zero except the combination that has all zero coefficients.

the market. Then we find the price of  $y$  as  $P = c_1P_1 + c_2P_2 + \dots + c_nP_n$ . Determining the correct linear combination can be a lot of work, however.

Remarkably, the CAPM in certainty equivalent form (that is, pricing form) does this work implicitly and correctly states that the price of any payoff  $y$  in the space  $\mathcal{M}$  generated by the  $y_i$ 's is given by

$$P_y = \frac{1}{R} [\bar{y} - \beta_{y,M} (\bar{M} - P_M R)], \quad (7.10)$$

where

$$\beta_{y,M} = \frac{\text{cov}(M, y)}{\sigma_M^2}$$

and where  $M$  is the Markowitz portfolio on the efficient frontier.

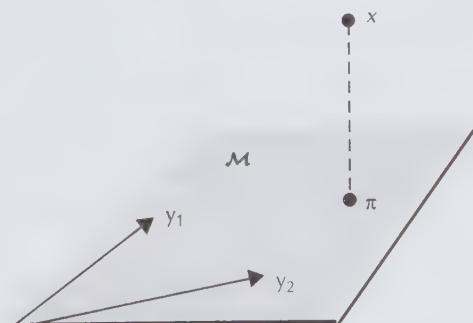
Suppose now that a new asset with payoff  $x$ , not in the space  $\mathcal{M}$ , is introduced. If we know the expected value and variance of  $x$  and its covariance with payoffs in  $\mathcal{M}$ , then  $x$  can be inserted into the right-hand side of equation (7.10) and the CAPM will produce a price. But, we may ask, how is this price related to the prices of payoffs in  $\mathcal{M}$ ? To characterize this price we introduce a simple geometric concept in the market space framework.

Imagine (in the two-dimensional case) the marketed space to be like the floor of a large room and the new asset payoff  $x$  to be a point above the floor. One way to price  $x$  is to find the point  $\pi$  on the floor that is closest to  $x$ . It will be straight down, so the line from  $x$  to this new point is perpendicular to the entire floor. That is,  $x - \pi$  is perpendicular (or orthogonal) to the floor. We can then assign the price to  $x$  as equal to  $\pi$ . See Figure 7.6.

To make this idea concrete we must define a few geometric concepts. The **inner product** of two random elements  $y_1$  and  $y_2$  is defined as  $E(y_1 y_2)$ . This is analogous to the dot product in ordinary vector geometry. We say two elements  $y_1$  and  $y_2$  are **orthogonal** if  $E(y_1 y_2) = 0$ . Also, for any element  $y$ , the quantity  $\sqrt{E(y)^2}$ , termed the **norm** of  $y$ , is denoted  $\|y\|$ . The norm serves as a measure of length. All of these notions can be applied to elements of  $\mathcal{M}$  and the space one dimension larger, which includes  $x$ .

A result known as the **projection theorem** states that, in general, there is a unique point  $\pi$  in  $\mathcal{M}$  closest to a given point  $x$ . Furthermore  $x - \pi$  is orthogonal to every point in  $\mathcal{M}$ .

**FIGURE 7.6** The element in  $\mathcal{M}$  closest to  $x$  is  $\pi$ , the projection of  $x$  onto  $\mathcal{M}$ . The projection always exists and is unique. Furthermore, the difference element  $x - \pi$ , is orthogonal to  $\mathcal{M}$ ; that is, for every payoff  $y$  in the market  $\mathcal{M}$ , there holds  $E[(x - \pi)y] = 0$ .



The point  $\pi$  is termed the **projection** of  $x$  onto  $\mathcal{M}$ . This provides us with the characterization of the price given by application of the CAPM. It is very easy to see that when pricing  $x$ , CAPM gives the price of its projection  $\pi$ .

**Projection price relation** *The CAPM price  $P_x$  of an asset with payoff  $x$  is the price of the projection  $\pi$  of  $x$  in the space of marketed asset payoffs.*

**Proof:** First note that  $\bar{x} = \bar{\pi}$ , because  $x - \pi$ , is orthogonal to the risk-free asset  $R$ . (That is,  $0 = E[(x - \pi)R] = (\bar{x} - \bar{\pi})R \Rightarrow \bar{x} = \bar{\pi}$ .) Similarly,  $\text{cov}(x, M) = \text{cov}(\pi, M)$ , because  $x - \pi$  is orthogonal to every element in  $\mathcal{M}$ . It follows (with  $M$  equal to the Markowitz efficient element in  $\mathcal{M}$ ) that

$$P_x \equiv \frac{1}{R}[\bar{x} - \beta_{x,M}(\bar{M} - P_M R)] = \frac{1}{R}[\bar{\pi} - \beta_{\pi,M}(\bar{M} - P_M R)] = P_\pi.$$

The first equality is the CAPM formula applied to  $x$  outside of  $\mathcal{M}$ . The second is when the formula is applied to  $\pi$ , which is in  $\mathcal{M}$ . ■

In view of this result, we say that CAPM pricing is equivalent to **projection pricing**.

Of course if  $x \in \mathcal{M}$ , then  $\pi = x$ . It is important to realize that the projection  $\pi$  always exists and is unique. By contrast, the usual CAPM (Markowitz) portfolio does not always exist, even though the projection price is always well defined. (See Exercise 16.)

## Minimum Norm Pricing\*

Now<sup>5</sup> that we have defined projection pricing, we can look for alternative methods that produce the projection price and that in some cases may be more convenient than the standard method. Again we use the concept of *orthogonality*.

Consider the set of assets in  $\mathcal{M}$  that have price 1. This will be linear surface. Define the element  $m_1$  as the payoff in this set of minimum norm; that is,  $m_1$  is closest to the zero payoff. Now, the projection theorem, viewed from the origin, says that the inner product  $E[m_1 y]$  will be constant for all  $y$ 's with price 1 because of the orthogonality relation shown<sup>6</sup> in Figure 7.7. We scale  $m_1$  to, say,  $M_1$  such that for some (or any)  $y_1$  with price 1 there holds  $E[M_1 y_1] = P_{y_1}$ .

The asset  $M_1$  then can be used to find the price of any asset in  $\mathcal{M}$ , as

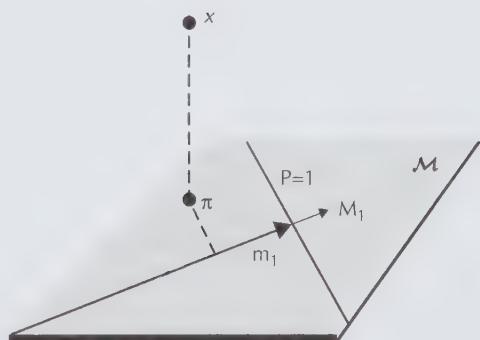
$$P_y = E[M_1 y]. \quad (7.11)$$

This is true because (7.11) correctly prices  $y_1$  and any other  $y_2$  with price 1. Any other asset  $y$  is simply a multiple of some asset of price 1, and since (7.11) is linear in  $y$ , it will price it correctly. Hence, this formula correctly prices everything in  $\mathcal{M}$ . We call  $M_1$  a **pricing element**. For linear pricing, there is always such an element, and we have shown one way that it can be constructed.

<sup>5</sup> This subsection can be omitted at first reading.

<sup>6</sup> If  $y_1$  and  $y_2$  have price 1, then  $E[m_1(y_1 - y_2)] = 0$ . So  $m_1$  is orthogonal to all assets that have price zero.

**FIGURE 7.7** The element  $m_1 \in \mathcal{M}$ . The inner product of this with all other elements of price 1 is constant. If  $m_1$  is scaled appropriately, to, say,  $M_1$ , its inner product with any element in  $\mathcal{M}$  will give the price of that element. That is,  $P_y = E[M_1 y]$  for all  $y$  in the market and will give the projection price for  $x$  outside the market.



Furthermore, for a new asset  $x$  outside the market  $\mathcal{M}$ , formula (7.11) gives the projection price of  $x$ . This is because  $x - \pi$  is orthogonal to  $\mathcal{M}$  and hence is orthogonal to  $M_1$ . That is,  $E[M_1(x - \pi)] = 0$ ; so  $E[M_1x] = E[M_1\pi]$ . Hence, even for  $x$  outside  $\mathcal{M}$ , we have the simple formula

$$P_x = E[M_1x] = P_\pi. \quad (7.12)$$

This is another way to calculate the projection price of an  $x$  that is not in the market space. When there is a risk-free asset with return  $R$ , it is possible to make (7.12) more *complicated*, expressing the pricing formula in the familiar form

$$P_x = \frac{1}{R} \left\{ \bar{x} - \frac{\text{cov}(x, M)}{\sigma_M^2} [\bar{M} - P_M R] \right\}, \quad (7.13)$$

where for simplicity we have omitted the subscript on  $M$ . Surprisingly, this is identical to (7.11). To verify this formula, see Exercise 12.

Notice that formula (7.13) looks exactly like the usual Markowitz/CAPM pricing formula in certainty equivalent form, except that the  $M$  in this result is a (scaled) minimum norm element rather than an efficient portfolio. Either one of these can be used. Both give the projection price.<sup>7</sup>

## 7.10 Correlation Pricing

When assigning a price to a new asset, say, a house, a person usually does not turn to a history of the S&P 500 stock index and attempt to correlate that index with the house in order to use the CAPM. Instead, one compares the house to other houses in the neighborhood that have recently been sold and then makes price adjustments to account for differences in size, location, and so forth. This **method of comparables** is used in everyday life for pricing everything from soap to TV sets; from used cars to private high-tech companies seeking to go public. Indeed, comparison is the most

<sup>7</sup> An advantage of more complex formula (7.13) over (7.12) is that (7.13) has scaling properties. In particular, if  $M_0$  is replaced by  $aM_0 + b$  in (7.13), the result does not change.

common pricing method. As we shall see, the method of comparison is one way to implement projection pricing and hence is a way to implement CAPM.

To price a new asset  $x$  by comparison, we look for marketed assets that are highly correlated with  $x$ . Ideally, we find an asset that, among all assets in the market space, is most correlated with  $x$ . This most correlated asset may be a single asset or a portfolio made up of a combination of assets. For a house, it may be a composite of houses in the same neighborhood. For a company, it may be a mixture of existing companies in the same industry. This most correlated asset (generally a portfolio) serves as a comparison, and it can replace the Markowitz portfolio in the pricing formula.

To proceed formally within our general framework, let us assert that a particular asset  $m \in \mathcal{M}$  is most correlated with  $x$ . Specifically,  $m \in \mathcal{M}$  is a solution to

$$\text{maximize } \text{cov}(x, m) \quad \text{subject to } \text{var}(m) \leq s^2,$$

where  $s > 0$  is an arbitrary positive constant that merely sets the scale for the correlated asset. Any multiple of the result is also most correlated with  $x$ .

We can describe the relation between a most correlated asset and the projection  $\pi$ . For any  $m \in \mathcal{M}$  with variance  $s^2$  we have

$$\begin{aligned} E[(x - \bar{x})(m - \bar{m})] &= E[(x - \pi + \pi - \bar{x})(m - \bar{m})] \\ &= E[(\pi - \bar{x})(m - \bar{m})] \\ &= E[(\pi - \bar{\pi})(m - \bar{m})] \\ &\leq ||\pi - \bar{\pi}|| \cdot ||m - \bar{m}|| = \sigma_\pi s. \end{aligned}$$

The second line follows from the orthogonality of  $x - \pi$  to  $\mathcal{M}$ . The third follows because  $\bar{\pi} = \bar{x}$ . The last line follows from the Cauchy–Schwarz inequality<sup>8</sup> and the fact that the optimal  $m$  will have standard deviation  $s$ .

One solution (with equality) is  $m = (s/\sigma_\pi)\pi$ . Any simple scaling of  $\pi$  is also most correlated with  $x$ . However, a constant (that is, risk-free value) can be added or subtracted, and the result is still a most correlated asset. Also, any positive multiple of a most correlated asset is also most correlated. Hence any most correlated asset is of the form  $M = a\pi + b$  for some  $a > 0$  and  $b$ . We now state the correlation pricing formula.

**Correlation Pricing Formula** *Let  $M$  be a marketed asset most correlated with  $x$ . Then the projection price of  $x$  may be computed as*

$$P_x = \frac{1}{R} \left\{ \bar{x} - \beta_{x,M} [\bar{M} - P_M R] \right\}, \quad (7.14)$$

---

<sup>8</sup> For any  $x$  and  $y$  there holds  $E[xy] \leq ||x|| \cdot ||y||$ .

where

$$\beta_{x,M} = \text{cov}(x, M) / \sigma_M^2.$$

**Proof:** Suppose first that  $M = \pi$ . Then, using  $\bar{x} = \bar{\pi}$  and  $\text{cov}(x, \pi) = \sigma_\pi^2$ , the formula gives

$$P_x \equiv \frac{1}{R} \{ \bar{x} - \beta_{x,M} [\bar{M} - P_M R] \} = \frac{1}{R} \{ \bar{\pi} - 1 \times [\bar{\pi} - P_\pi R] \} = P_\pi,$$

as it should. For  $M$  of the form  $a\pi + b$ , the same result will be obtained because of the invariance of the formula to such transformations. See Exercise 14. ■

In conclusion we now have *three* logical choices for  $M$ , all of which are abbreviated by the letter “M”: (1) The original Markowitz efficient portfolio, which is used in CAPM; (2) the minimum-norm portfolio, which always exists; and (3) an asset most correlated with the payoff  $x$  we wish to price. Choices (1) and (2) are universal, in the sense that each of them applies to any new payoff, while (3) depends on  $x$ . All three choices give the same result. In practice, some sort of comparison procedure is most commonly used for pricing.

**Example 7.8 (In the market)** Suppose that  $x$  happens to be in the market. Then an asset most correlated with  $x$  is  $x$  itself. Hence, we may substitute  $x$  for  $M$  in the pricing formula. This gives

$$P_x = \frac{1}{R} \{ \bar{x} - 1 \times [\bar{x} - P_x R] \} = \frac{1}{R} \{ \bar{x} - \bar{x} + P_x R \} = P_x,$$

as expected.

**Example 7.9 (New software)** The management team of a software firm is considering the development of a new product, and wants to assign a suitable value to it. A one-year perspective is appropriate. The net revenue  $x$  from the product (in one year) has expected value \$100,000. The current rate of interest is 5%. The team feels that to find a present value they should discount more heavily, so they meet to decide on an appropriate rate. Initially they settle on 17%, since that seems to be what many other companies use. However, suspicious of simple discounting methods, the CEO diplomatically suggests that the value should be the projection price of the project. After quick agreement, the team realizes that a most-correlated marketed security may be the value of the firm's own stock. Historical analysis suggests that the correlation coefficient between such a product and the stock price is about  $\rho = 0.8$ . The current value of the stock is  $M = \$100$  per share and it is believed that this will increase by 15% over the coming year. The ratio of standard deviations  $\sigma_x / \sigma_M$  is believed to be 900, which is slightly higher than  $\bar{x}/\bar{M} = 100,000/115 = 869$  to reflect the project's

greater proportionate risk. The value of the project is therefore found to be

$$\begin{aligned}
 P_x &= \frac{1}{R} \left\{ \bar{x} - \frac{\sigma_{xM}}{\sigma_M^2} [\bar{M} - P_M R] \right\} \\
 &= \frac{1}{1.05} \left\{ 100,000 - \frac{\sigma_{xM}}{\sigma_x \sigma_M} \frac{\sigma_x}{\sigma_M} [115 - 105] \right\} \\
 &= \frac{1}{1.05} \left\{ 100,000 - 0.8 \cdot 900 \cdot 10 \right\} \\
 &= \frac{1}{1.05} \left\{ 100,000 - 7,200 \right\} = 92,800 / 1.05 = \$88,381.
 \end{aligned}$$

This amounts to a discount of expected value by 13.1%.

## 7.11 Summary

If everybody uses the mean–variance approach to investing, and if everybody has the same estimates of the asset's expected returns, variances, and covariances, then everybody must invest in the same fund  $F$  of risky assets and in the risk-free asset. Because  $F$  is the same for everybody, it follows that, in equilibrium,  $F$  must correspond to the market portfolio  $M$ —the portfolio in which each asset is weighted by its proportion of total market capitalization. This observation is the basis for the capital asset pricing model (CAPM).

If the market portfolio  $M$  is the efficient portfolio of risky assets, it follows that the efficient frontier in the  $\bar{r}$ – $\sigma$  diagram is a straight line that emanates from the risk-free point and passes through the point representing  $M$ . This line is the capital market line. Its slope is called the market price of risk. Any efficient portfolio must lie on this line.

The CAPM is derived directly from the condition that the market portfolio is a point on the edge of the feasible region that is tangent to the capital market line; in other words, the CAPM expresses the tangency conditions in mathematical form. The CAPM result states that the expected rate of return of any asset  $i$  satisfies

$$\bar{r}_i - r_f = \beta_i (\bar{r}_M - r_f),$$

where  $\beta_i = \text{cov}(r_i, r_M) / \sigma_M^2$  is the beta of the asset.

The CAPM can be represented graphically as a security market line: the expected rate of return of an asset is a straight-line function of its beta (or, alternatively, of its covariance with the market); greater beta implies greater expected return. Indeed, from the CAPM view it follows that the risk of an asset is fully characterized by its beta. It follows, for example, that an asset that is uncorrelated with the market ( $\beta = 0$ ) will have an expected rate of return equal to the risk-free rate.

The beta of the market portfolio is by definition equal to 1. The betas of other stocks take other values, but the betas of most U.S. stocks range between .5 and 2.5. The beta of a portfolio of stocks is equal to the weighted average of the betas of the individual assets that make up the portfolio.

One application of CAPM is the evaluation of mutual fund performance. The Jensen index measures the historical deviation of a fund from the security market line. (This measure has dubious value for funds of publicly traded stocks, however.) The Sharpe index measures the slope of the line joining the fund and the risk-free asset on the  $\bar{r}-\sigma$  diagram, so that this slope can be compared with the market price of risk.

The CAPM can be converted to an explicit formula for the price of an asset. In the simplest version, this formula states that price is obtained by discounting the expected payoff, but the interest rate used for discounting must be  $r_f + \beta(\bar{r}_M - r_f)$ , where  $\beta$  is the beta of the asset. An alternative form expresses the price as a discounting of the certainty equivalent of the payoff, and in this formula the discounting is based on the risk-free rate  $r_f$ .

It is important to recognize that the pricing formula of CAPM is linear, meaning that the price of a sum of assets is the sum of their prices, and the price of a multiple of an asset is that same multiple of the basic price. The certainty equivalent formulation of the CAPM clearly exhibits this linear property.

The CAPM is frequently applied to new assets that are not yet part of the market. In this circumstance the CAPM formula produces a price equal to the price of an asset in the market whose payoff is closest to the payoff of the new asset in the sense of minimum norm (with the norm of a payoff  $y$  being  $\|y\| = \sqrt{E[y^2]}$ ). If the new payoff is  $x$  and the closest marketed payoff is  $\pi$ , then the error payoff  $x - \pi$  satisfies  $E[(x - \pi)m] = 0$  for all market payoffs  $m$ —a property stated geometrically as saying that the error is orthogonal to the market. This pricing is called **projection pricing**, since the closest payoff in the market is, geometrically, the projection of the  $x$  onto the collection of market payoffs.

There are other ways to find the projection price. One way is to use the correlation pricing formula. In this method, a portfolio  $M$  in the market most correlated with a payoff  $x$  is used in the standard formula to find the projection price of  $x$ . This method provides a rigorous justification for pricing a new asset by comparing it with existing assets of a similar nature.

## Exercises

1. (Capital market line) Assume that the expected rate of return on the market portfolio is 23% and the rate of return on T-bills (the risk-free rate) is 7%. The standard deviation of the market is 32%. Assume that the market portfolio is efficient.
  - (a) What is the equation of the capital market line?
  - (b) (i) If an expected return of 39% is desired, what is the standard deviation of this position? (ii) If you have \$1,000 to invest, how should you allocate it to achieve the above position?
  - (c) If you invest \$300 in the risk-free asset and \$700 in the market portfolio, how much money should you expect to have at the end of the year?
2. (A small world) Consider a world in which there are only two risky assets,  $A$  and  $B$ , and a risk-free asset  $F$ . The two risky assets are in equal supply in the market; that is,  $M = \frac{1}{2}(A + B)$ . The following information is known:  $r_F = 10$ ,  $\sigma_A^2 = 04$ ,  $\sigma_{AB} = 01$ ,  $\sigma_B^2 = .02$ , and  $\bar{r}_M = .18$ .

- (a) Find a general expression (without substituting values) for  $\sigma_M^2$ ,  $\beta_A$ , and  $\beta_B$ .  
 (b) According to the CAPM, what are the numerical values of  $\bar{r}_A$  and  $\bar{r}_B$ ?

3. (Bounds on returns) Consider a universe of just three securities. They have expected rates of return of 10%, 20%, and 10%, respectively. Two portfolios are known to lie on the minimum-variance set. They are defined by the portfolio weights

$$\mathbf{w} = \begin{bmatrix} .60 \\ .20 \\ .20 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} .80 \\ -.20 \\ .40 \end{bmatrix}$$

It is also known that the market portfolio is efficient.

- (a) Given this information, what are the minimum and maximum possible values for the expected rate of return on the market portfolio?  
 (b) Now suppose you are told that  $\mathbf{w}$  represents the minimum-variance portfolio. Does this change your answers to part (a)?

4. (Quick CAPM derivation) Derive the CAPM formula for  $\bar{r}_k - r_f$  by using Equation (6.9) in Chapter 6. [Hint: Note that

$$\sum_{i=1}^n \sigma_{ik} w_i = \text{cov}(r_k, r_M).$$

Apply equation (6.9) both to asset  $k$  and to the market itself.]

5. (Uncorrelated assets) Suppose there are  $n$  mutually uncorrelated assets. The return on asset  $i$  has variance  $\sigma_i^2$ . The expected rates of return are unspecified at this point. The total amount of asset  $i$  in the market is  $X_i$ . We let  $T = \sum_{i=1}^n X_i$  and then set  $x_i = X_i/T$ , for  $i = 1, 2, \dots, n$ . Hence the market portfolio in normalized form is  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ . Assume there is a risk-free asset with rate of return  $r_f$ . Find an expression for  $\beta_j$  in terms of the  $x_i$ 's and  $\sigma_i$ 's.

6. (Simpleland) In Simpleland there are only two risky stocks, A and B, whose details are listed in Table 7.4.

**TABLE 7.4**  
**DETAILS OF STOCKS A AND B**

	Number of shares outstanding	Price per share	Expected rate of return	Standard deviation of return
Stock A	100	\$1.50	15%	15%
Stock B	150	\$2.00	12%	9%

Furthermore, the correlation coefficient between the returns of stocks A and B is  $\rho_{AB} = \frac{1}{3}$ . There is also a risk-free asset, and Simpleland satisfies the CAPM exactly.

- (a) What is the expected rate of return of the market portfolio?  
 (b) What is the standard deviation of the market portfolio?  
 (c) What is the beta of stock A?  
 (d) What is the risk-free rate in Simpleland?

7. (Zero-beta assets) Let  $\mathbf{w}_0$  be the portfolio (weights) of risky assets corresponding to the minimum-variance point in the feasible region. Let  $\mathbf{w}_1$  be any other portfolio on the efficient frontier. Define  $r_0$  and  $r_1$  to be the corresponding returns.

- (a) There is a formula of the form  $\sigma_{01} = A\sigma_0^2$ . Find  $A$ . [Hint: Consider the portfolios  $(1 - \alpha)\mathbf{w}_0 + \alpha\mathbf{w}_1$ , and consider small variations of the variance of such portfolios near  $\alpha = 0$ .]
- (b) Corresponding to the portfolio  $\mathbf{w}_1$  there is a portfolio  $\mathbf{w}_z$  on the minimum-variance set that has zero beta with respect to  $\mathbf{w}_1$ ; that is,  $\sigma_{1,z} = 0$ . This portfolio can be expressed as  $\mathbf{w}_z = (1 - \alpha)\mathbf{w}_0 + \alpha\mathbf{w}_1$ . Find the proper value of  $\alpha$ .
- (c) Show the relation of the three portfolios on a diagram that includes the feasible region.
- (d) If there is no risk-free asset, it can be shown that other assets can be priced according to the formula

$$\bar{r}_i - \bar{r}_z = \beta_{iM}(\bar{r}_M - \bar{r}_z),$$

where the subscript  $M$  denotes the market portfolio and  $\bar{r}_z$  is the expected rate of return on the portfolio that has zero beta with the market portfolio. Suppose that the expected returns on the market and the zero-beta portfolio are 15% and 9%, respectively. Suppose that a stock  $i$  has a correlation coefficient with the market of 0.5. Assume also that the standard deviation of the returns of the market and stock  $i$  are 15% and 5%, respectively. Find the expected return of stock  $i$ .

8. (Wizards ◊) Electron Wizards, Inc. (EWI) has a new idea for producing TV sets, and it is planning to enter the development stage. Once the product is developed (which will be at the end of 1 year), the company expects to sell its new process for a price  $p$ , with expected value  $\bar{p} = \$24M$ . However, this sale price will depend on the market for TV sets at the time. By examining the stock histories of various TV companies, it is determined that the final sales price  $p$  is correlated with the market return as  $E[(p - \bar{p})(r_M - \bar{r}_M)] = \$20M\sigma_M^2$ .

To develop the process, EWI must invest in a research and development project. The cost  $c$  of this project will be known shortly after the project is begun (when a technical uncertainty will be resolved). The current estimate is that the cost will be either  $c = \$20M$  or  $c = \$16M$ , and each of these is equally likely. (This uncertainty is uncorrelated with the final price and is also uncorrelated with the market.) Assume that the risk-free rate is  $r_f = 9\%$  and the expected return on the market is  $\bar{r}_M = 33\%$ .

- (a) What is the expected rate of return of this project?
- (b) What is the beta of this project? [Hint: In this case, note that

$$E\left[\left(\frac{p - \bar{p}}{c}\right)(r_M - \bar{r}_M)\right] = E\left(\frac{1}{c}\right)E[(p - \bar{p})(r_M - \bar{r}_M)].$$

- (c) Is this an acceptable project based on a CAPM criterion? In particular, what is the excess rate of return (+ or -) above the return predicted by the CAPM?

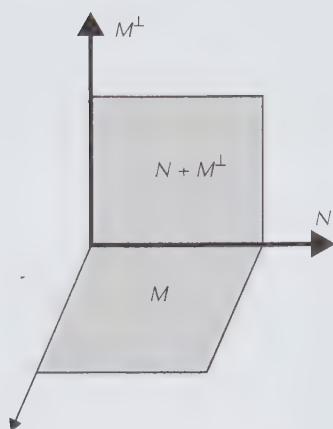
9. (Gavin's problem) Prove to Gavin Jones that the results he obtained in Examples 7.5 and 7.7 were not accidents. Specifically, for a fund with return  $\alpha r_f + (1 - \alpha)r_M$ , show that both CAPM pricing formulas give the price of \$100 worth of fund assets as \$100.

10. (Little world CAPM) Suppose there are only two risky assets with expected rates of return  $\bar{r}_1 = 0.1, \bar{r}_2 = 0.2$  and covariances  $\sigma_1^2 = 0.04, \sigma_2^2 = 0.09, \sigma_{1,2} = 0.03$ . The current risk-free rate is 0.05.

- (a) Find the one fund for this small market.

- (b) Now suppose that there is a new project with an expected payoff  $\bar{Q}_1 = 10$  and covariances  $\sigma_Q^2 = 25, \sigma_{Q,1} = 0.8, \sigma_{Q,2} = 0.15$ . Using the one fund as the market portfolio, what is the price of this project according to CAPM?
11. (Risk analysis) Assume the market portfolio has expected rate of return  $\bar{r}_m = 0.12$  and standard deviation  $\sigma_m = 0.3$ . The risk-free rate is  $r_f = 0.02$ . There is another stock, a, in the market with  $\bar{\sigma}_a = 0.6, \rho_{am} = 0.1$ .
- Find  $\bar{r}_a$  and  $\beta_a$ .
  - A new asset, b, has the same expected return as a but a standard deviation of  $\sigma_b = 0.8$ . What is the idiosyncratic error of b?
  - Another asset, c, enters the market with  $\rho_c = 0.8$ . What percentage of the risk of c is idiosyncratic?
12. (Standard form) Let  $M$  be a marketed asset that is also a pricing asset such that for every marketed payoff  $x$  there holds  $P_x = E[Mx]$ . Show that it follows that
- $$P_x = \frac{1}{R} \left\{ \bar{x} - \frac{\text{cov}(x, M)}{\sigma_M^2} [\bar{M} - P_M R] \right\}.$$
- [Hint: Apply the original formula to 1 and to  $M$  itself.]
13. (A leveraged firm) A company earns a rate of return of  $r_A$  and has beta  $\beta_A$ . A fraction  $w$  of the assets is owned by bondholders, and the remaining fraction  $(1-w)$  is owned by equity holders. Every year, the bondholders demand a riskless rate of return of  $r_B$  on their fraction of the assets, regardless of the actual rate of return  $r_A$  that was achieved that year. Beyond that, the equity holders take whatever is left after the bondholders have been paid.
- What is the rate of return of the equity holders in terms of  $w, r_A$ , and  $r_B$ ?
  - What is the beta of the rate of return of the equity holders in terms of  $w$  and  $\beta_A$ ?
  - Suppose  $\beta_A$  is positive and the expected rate of return on the market is greater than the risk-free rate. As  $w$  increases (that is, as the firm becomes more leveraged), what should happen to the expected rate of return on the equity in the firm?
14. (Equivalent forms) Consider the pricing formula
- $$P_x = \frac{1}{R} \left\{ \bar{x} - \beta_{x,M} (\bar{M} - P_M R) \right\}.$$
- Let  $N = aM + b$  for constants  $a > 0$  and  $b$ . Show that substitution of  $N$  for  $M$  yields the exact same formula but with  $N$  replacing  $M$ .
15. (Singular) Suppose there are two stocks that are uncorrelated. Each of these has variance of 1, and there are expected returns are  $\bar{r}_1$  and  $\bar{r}_2$ , respectively. The risk-free rate is  $r_f$ . Find the portfolio of weights  $w_1$  and  $w_2$  for the Markowitz (market) portfolio. Show that for some value of  $r_f$  there is no Markowitz portfolio.
16. (Separate spaces\*) Suppose there is a master space  $\Omega$  of payoff elements (of finite dimension), of which those in the market are in the subspace  $\mathcal{M}$ .
- We may write  $\Omega = \mathcal{M} + \mathcal{M}^\perp$ , where  $\mathcal{M}^\perp$  is the set of all elements orthogonal to  $\mathcal{M}$ . That is, any element  $x \in \Omega$  can be expressed as  $x = m + m^\perp$  with  $m \in \mathcal{M}$  and  $m^\perp \in \mathcal{M}^\perp$ . Let  $\mathcal{N}$  be a subspace of  $\mathcal{M}$ , and consider all payoffs in the subspace  $\mathcal{X} = \mathcal{N} + \mathcal{M}^\perp$  in

**FIGURE 7.8** The subspace  $\mathcal{M}$  is a whole plane of which just a square segment is shown.  $\mathcal{M}^\perp$  is the vertical axis, and  $\mathcal{N}$  is the horizontal axis.



$\Omega$ . (See Figure 7.8 for  $\Omega$  equal to three-dimensional space.) Show that for an  $x \in \mathcal{X}$ , the projection of  $x$  onto  $\mathcal{M}$  is equal to the projection of  $x$  onto  $\mathcal{N}$ . Likewise, a most correlated asset to  $x$  is in  $\mathcal{N}$ . Hence, to find the projection price of  $x$ , it is only necessary to consider its relation to market payoffs in its local  $\mathcal{N}$ .

17. (Two assets) Suppose there are two marketed assets, each with price 1 and payoffs  $y_1$  and  $y_2$ , respectively, with  $\bar{y}_1 = 1.4$  and  $\bar{y}_2 = 0.8$ . Each has a variance of 0.04, and they are uncorrelated.
  - (a) Find the minimum-norm portfolio.
  - (b) Find the projection price of the risk-free asset with payoff 1.
18. (Going public) Firm X is about to be offered publicly, and the investment banking firm that is facilitating the offering is working to determine an appropriate offering price. There is a publicly traded company Y that is quite similar to X. The 1-year value of Y (the total value of all shares of stock) is expected to be \$500 million with a volatility of 20%, and that of X is expected to be \$100 million with a volatility of 30%. The coefficient of correlation between X and Y is  $\rho = 0.8$ . The interest rate is 10%, and the growth rate of Y's price is expected to be 20%. That is, currently  $P_Y G = \bar{Y}$ , where  $G$  is the expected total return equal to 1 plus the expected growth rate. What is the appropriate price for firm X, and what is the discount rate that is equivalent to this price?

## References

The CAPM theory was developed independently in references [1-4]. There are now numerous extensions and textbook accounts of that theory. Consult any of the basic finance textbooks listed as references for Chapter 2. The application of this theory to mutual fund performance evaluation was presented in [5, 6]. An alternative measure, not discussed in this chapter, is due to Treynor [7]. For summaries of the application of CAPM to corporate analysis, see [8, 9]. The idea of using a zero-beta asset, as in Exercise 7, is due to Black [10]. Projection theory is treated in detail in [11]. Use of the minimum-norm portfolio for finding the projection price is discussed in [12], where the example of Exercise 12 is also presented. The correlation pricing formula is formulated in [13].

1. Sharpe, W. F. (1964), "Capital Asset Prices: A Theory of Market Equilibrium under Conditions of Risk," *Journal of Finance*, **19**, 425–442.
2. Lintner, J. (1965), "The Valuation of Risk Assets and the Selection of Risky Investment in Stock Portfolios and Capital Budgets," *Review of Economics and Statistics*, **47**, 13–37.
3. Mossin, J. (1966), "Equilibrium in a Capital Asset Market," *Econometrica*, **34**, no. 4, 768–783.
4. Treynor, J. L. (1961), "Towards a Theory of Market Value of Risky Assets," unpublished manuscript.
5. Sharpe, W. F. (1966), "Mutual Fund Performance," *Journal of Business*, **39**, January, 119–138.
6. Jensen, M. C. (1969), "Risk, the Pricing of Capital Assets, and the Evaluation of Investment Portfolios," *Journal of Business*, **42**, April, 167–247.
7. Treynor, J. L. (1965), "How to Rate Management Investment Funds," *Harvard Business Review*, **43**, January–February, 63–75.
8. Rubinstein, M. E. (1973), "A Mean–Variance Synthesis of Corporate Financial Theory," *Journal of Finance*, **28**, 167–182.
9. Fama, E. F. (1977), "Risk-Adjusted Discount Rates and Capital Budgeting under Uncertainty," *Journal of Financial Economics*, **5**, 3–24.
10. Black, F. (1972), "Capital Market Equilibrium with Restricted Borrowing," *Journal of Business*, **45**, 445–454.
11. Luenberger, David G. (1969), *Optimization by Vector Space Methods*, Wiley, New York.
12. Luenberger, David G. (2002), "Projection Pricing," *Journal of Optimization Theory and Applications*, **109**, 1–25.
13. Luenberger, David G. (2002), "A Correlation Pricing Formula," *Journal of Economic Dynamics and Control*, **26**, 1113–1126.

# 8

## OTHER PRICING MODELS

### 8.1 Introduction

The theory of the previous two chapters is quite general, for it can be applied to bets on a wheel of fortune, to analysis of an oil wildcat venture, to construction of a portfolio of stocks, and to many other single-period investment problems. However, the primary application of mean–variance theory is to stocks, and this chapter focuses primarily on those special securities, although much of the material is applicable to other assets as well.

This chapter examines how models of stock returns, suitable for mean–variance analysis, can be specified. It shows how to build a **factor model** of the return process to simplify the structure and reduce the number of required parameters. Along the way a new theory of asset pricing, termed **arbitrage pricing theory** (APT), is obtained. In the next chapter, we turn directly to the issue of determining parameter values. We consider the possibility of using historical data to determine parameter values, but we discover that this approach is of limited value.

### 8.2 Factor Models

The information required by the mean–variance approach grows substantially as the number  $n$  of assets increases. There are  $n$  mean values,  $n$  variances, and  $n(n - 1)/2$  covariances—a total of  $2n + n(n - 1)/2$  parameters. When  $n$  is large, this is a very large set of required values. For example, if we consider a universe of 1,000 stocks,

501,500 values are required to fully specify a mean–variance model. Clearly it is a formidable task to obtain this information directly. We need a simplified approach.

Fortunately the randomness displayed by the returns of  $n$  assets often can be traced back to a smaller number of underlying basic sources of randomness (termed factors) that influence the individual returns. A factor model that represents this connection between factors and individual returns leads to a simplified structure for the covariance matrix, and provides important insight into the relationships among assets.

The factors used to explain randomness must be chosen carefully—and the proper choice depends on the universe of assets being considered. For real estate parcels within a city, the underlying factors might be population, employment rate, and school budgets. For common stocks listed on an exchange, the factors might be the stock market average, gross national product, employment rate, and so forth. Selection of factors is somewhat of an art, or a trial-and-error process, although formal analysis methods can also be helpful. (See Exercise 3.)

This section introduces the factor model concept and shows how it simplifies the covariance structure.

## Single-Factor Model

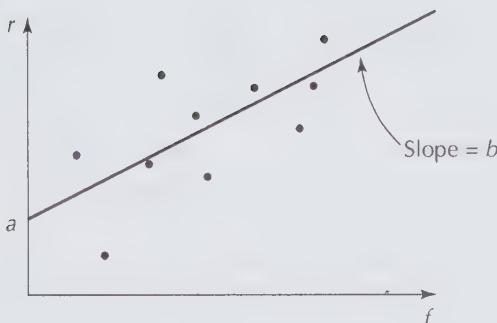
Single-factor models are the simplest of the factor models, but they illustrate the concept quite well. Suppose that there are  $n$  assets, indexed by  $i$ , with rates of return  $r_i, i = 1, 2, \dots, n$ . There is a single factor  $f$  which is a random quantity (such as the stock market average rate of return for the period). We assume that the rates of return and the factor are related by the following equation:

$$r_i = a_i + b_i f + e_i \quad (8.1)$$

for  $i = 1, 2, \dots, n$ . In this equation, the  $a_i$ 's and the  $b_i$ 's are fixed constants. The  $e_i$ 's are random quantities which represent **errors**. Without loss of generality, it can be assumed that the errors each have zero mean, that is,  $E(e_i) = 0$ , since any nonzero mean could be transferred to  $a_i$ . In addition, however, it is usually assumed that the errors are uncorrelated with  $f$  and with each other; that is,  $E[(f - \bar{f})e_i] = 0$  for each  $i$  and  $E(e_i e_j) = 0$  for  $i \neq j$ . These are idealizing assumptions which may not actually be true, but are usually assumed to be true for purposes of analysis. It is also assumed that variances of the  $e_i$ 's are known, and they are denoted by  $\sigma_{e_i}^2$ .

An individual factor model equation can be viewed graphically as defining a linear fit to (potential) data, as shown in Figure 8.1. Imagine that several independent observations are made of both the rate of return  $r_i$  and the factor  $f$ . These points are plotted on the graph. Since both are random quantities, the points are likely to be scattered. A straight line is fitted through these points in such a way that the average value of the error, as measured by the vertical distance from a point to the line, is zero.

It is helpful to view Figure 8.1 in two ways. First, given the model (8.1), we can draw the line on the diagram before obtaining data points. Then if we believe the model, we believe that the data points will fall in the kind of pattern shown in the diagram. In the second view, we imagine that we first obtain the data points, then we construct the line that fits the data. When we draw the line, however, we are implying that additional data are likely to support it in the sense of falling in the same pattern.



**FIGURE 8.1 Single-factor model.** Returns are related linearly to the factor  $f$ , except that random errors are added to the return.

When applied to a group of assets, the fitting process is carried out for each asset separately. As a result, we obtain for each asset  $i$  an  $a_i$  and  $b_i$ . The  $a_i$ 's are termed **intercepts** because  $a_i$  is the intercept of the line for asset  $i$  with the vertical axis. The  $b_i$ 's are termed **factor loadings** because they measure the sensitivity of the return to the factor.

If an historical record of asset returns and the factor values are available, the parameters of a single-factor model can be estimated by actually fitting straight lines, as suggested before. Note, however, that different values of the  $a_i$ 's and  $b_i$ 's are likely to be obtained for different sets of data. For example, if we use monthly data on returns and the factor  $f$  for one year to obtain values of the  $a_i$ 's and  $b_i$ 's, and then we do it again the next year, we are likely to get different values. In what follows, we assume that the model is given, and that it represents our understanding of how the returns are related to the factor  $f$ . We ignore the question of where this model comes from—at least for now.

If we agree to use a single-factor model, then the standard parameters for mean-variance analysis can be determined directly from that model. We calculate

$$\bar{r}_i = a_i + b_i \bar{f} \quad (8.2a)$$

$$\sigma_i^2 = b_i^2 \sigma_f^2 + \sigma_{e_i}^2 \quad (8.2b)$$

$$\sigma_{ij} = b_i b_j \sigma_f^2, \quad i \neq j \quad (8.2c)$$

$$b_i = \text{cov}(r_i, f) / \sigma_f^2. \quad (8.2d)$$

These equations reveal the primary advantage of a factor model. In the usual representation of asset returns, a total of  $2n + n(n - 1)/2$  parameters are required to specify means, variances, and covariances. In a single-factor model, only the  $a_i$ 's,  $b_i$ 's,  $\sigma_{e_i}^2$ 's, and  $\bar{f}$  and  $\sigma_f^2$  are required—a total of just  $3n + 2$  parameters.

## Portfolio Parameters

When asset returns are described by a single-factor model, the return of any portfolio of these assets is described by a corresponding factor model equation of its own. To

verify this important property, suppose that there are  $n$  assets with rates of return governed by the factor model

$$r_i = a_i + b_i f + e_i, \quad i = 1, 2, \dots, n.$$

Suppose that a portfolio is constructed with weights  $w_i$ , with  $\sum_{i=1}^n w_i = 1$ . Then the rate of return  $r$  of the portfolio is just the corresponding combination of individual rates of return; namely,

$$r = \sum_{i=1}^n w_i a_i + \sum_{i=1}^n w_i b_i f + \sum_{i=1}^n w_i e_i.$$

We can write this as

$$r = a + b f + e,$$

where

$$a = \sum_{i=1}^n w_i a_i$$

$$b = \sum_{i=1}^n w_i b_i$$

$$e = \sum_{i=1}^n w_i e_i.$$

Both  $a$  and  $b$  are constants, which are weighted averages of the individual  $a_i$ 's and  $b_i$ 's. The error term  $e$  is random, but it, too, is an average. Under the assumptions that  $E(e_i) = 0$ ,  $E[(f - \bar{f})e_i] = 0$ , and  $E(e_i e_j) = 0$  for all  $i \neq j$ , it is clear that  $E(e) = 0$  and  $E[(f - \bar{f})e] = 0$ ; that is,  $e$  and  $f$  are uncorrelated. The variance of  $e$  is

$$\sigma_e^2 = E(e^2) = E\left[\left(\sum_{i=1}^n w_i e_i\right)\left(\sum_{j=1}^n w_j e_j\right)\right] = E\left(\sum_{i=1}^n w_i^2 e_i^2\right) = \sum_{i=1}^n w_i^2 \sigma_{e_i}^2$$

where we have used the fact that the  $e_i$ 's are uncorrelated with each other. Thus we have a simple and full description of the portfolio return as a factor equation.

A factor model is a good model to use to explore the effects of diversification, showing how risk can be reduced but not entirely eliminated. For simplicity, let us assume that in the one-factor model  $\sigma_{e_i}^2$  is the same for all  $i$ ; say,  $\sigma_{e_i}^2 = s^2$ . Suppose that a portfolio is formed by taking equal fractions of each asset; that is, we put  $w_i = 1/n$  for each  $i$ . In that case, from before, we find

$$\sigma_e^2 = \frac{1}{n} s^2.$$

Hence as  $n \rightarrow \infty$  we see that  $\sigma_e^2 \rightarrow 0$ . So in a well-diversified portfolio the error term in the factor equation is small.

The overall variance of the portfolio is

$$\sigma^2 = b^2 \sigma_f^2 + \sigma_e^2.$$

The  $\sigma_e^2$  term goes to zero, but since  $b$  is an average of the  $b_i$ 's, the  $b^2\sigma_f^2$  term remains more or less constant. Hence the variance of the portfolio tends to decrease as  $n$  increases because  $\sigma_e^2$  goes to zero, but the portfolio variance does not go to zero.

This observation leads to a general conclusion. For any one asset with a rate of return described by a factor model

$$r_i = a_i + b_i f + e_i$$

there are two sources of risk: that due to the  $b_i f$  term and that due to  $e_i$ . The risk due to  $e_i$  is said to be **diversifiable** because this term's contribution to overall risk is essentially zero in a well-diversified portfolio. On the other hand, the  $b_i f$  term is said to be a **systematic** or **nondiversifiable risk**, since it is present even in a diversified portfolio. The systematic risk is due to the factor that influences every asset, so diversification cannot eliminate it. The risks due to the  $e_i$ 's are independent and, hence, each can be reduced by diversification.

**Example 8.1 (Four stocks and one index)** The upper portion of Table 8.1 shows the historical rates of return (in percent) for four stocks over a period of 10 years. Also shown is a record of an industrial price index over this same period. We shall build a single-index model for each of the stocks using this index as the factor. As a first step, we calculate the historical averages of the returns and the index. We

**TABLE 8.1  
FACTOR MODEL**

Year	Stock 1	Stock 2	Stock 3	Stock 4	Index
1	11.91	29.59	23.27	27.24	12.30
2	18.37	15.25	19.47	17.05	5.50
3	3.64	3.53	-6.58	10.20	4.30
4	24.37	17.67	15.08	20.26	6.70
5	30.42	12.74	16.24	19.84	9.70
6	-1.45	-2.56	-15.05	1.51	8.30
7	20.11	25.46	17.80	12.24	5.60
8	9.28	6.92	18.82	16.12	5.70
9	17.63	9.73	3.05	22.93	5.70
10	15.71	25.09	16.94	3.49	3.60
aver	15.00	14.34	10.90	15.09	6.74
var	90.28	107.24	162.19	68.27	6.99
cov	2.34	4.99	5.45	11.13	6.99
$b$	0.33	0.71	0.78	1.59	1.00
$a$	12.74	9.53	5.65	4.36	0.00
e-var	89.49	103.68	157.95	50.55	

The record of the rates of return for four stocks and an index of industrial prices are shown. The averages and variances are all computed, as well as the covariance of each with the index. From these quantities, the  $b_i$ 's and the  $a_i$ 's are calculated. Finally, the computed error variances are also shown. The index does not explain the stock price variations very well.

denote the averages by  $\hat{r}_i$  and  $\hat{f}$  to distinguish these values from the true (but unknown) values  $\bar{r}_i$  and  $\bar{f}$ .

Let  $r_i^k$ , for  $k = 1, 2, \dots, 10$ , denote the 10 samples of the rate of return  $r_i$ . Then the estimate of  $\bar{r}_i$  is

$$\hat{r}_i = \frac{1}{10} \sum_{k=1}^{10} r_i^k.$$

We estimate the variances with the formula

$$\text{var}(r_i) = \frac{1}{9} \sum_{k=1}^{10} (r_i^k - \hat{r}_i)^2,$$

which is the standard way to estimate variance.<sup>1</sup> Analogous formulas are used to calculate estimates of the mean and the variance of the index.

Next the covariances of the returns with the index are estimated. The formula used for this purpose is

$$\text{cov}(r_i, f) = \frac{1}{9} \sum_{k=1}^{10} (r_i^k - \hat{r}_i)(f^k - \hat{f}). \quad (8.3)$$

Once the covariances are estimated, we find the values of  $b_i$  and  $a_i$  from the formulas

$$b_i = \frac{\text{cov}(r_i, f)}{\text{var}(f)}$$

$$a_i = \hat{r}_i - b_i \hat{f}.$$

(The first of these is obtained by forming the covariance with respect to  $f$  of both sides of the factor equation.)

After the model is constructed, we estimate the variance of the error under the assumption that the errors are uncorrelated with each other and with the index. Hence, using equation (8.2b) we write

$$\text{var}(e_i) = \text{var}(r_i) - b_i^2 \text{var}(f).$$

These values are shown in the last row of Table 8.1. Notice that these error variances are almost as large as the variances of the stock returns themselves, and hence the factor does not explain much of the variation in returns. In other words, there is high nonsystematic risk. Furthermore, by applying a version of (8.3) to estimate  $\text{cov}(e_i, e_j)$ , it turns out that the errors are highly correlated. For example, the estimation formula gives  $\text{cov}(e_1, e_2) = 44$  and  $\text{cov}(e_2, e_3) = 91$ , whereas the factor model was constructed under the assumption that these error covariances are zero. Hence this single-index model is not a very accurate representation of the stock returns. (A better model for these data is given in the next section.)

---

<sup>1</sup> See Section 9.2 for details on this estimation formula.

## Multifactor Models

The preceding development can be extended to include more than one factor. For example, if there are two factors  $f_1$  and  $f_2$ , with perhaps the first factor being a broad index of the market return and the second an index of the change since the previous period of consumer spending, the model for the rate of return of asset  $i$  would have the form

$$r_i = a_i + b_{1i}f_1 + b_{2i}f_2 + e_i.$$

Again the constant  $a_i$  is called the intercept, and  $b_{1i}$  and  $b_{2i}$  are the factor loadings. The factors  $f_1$  and  $f_2$  and the error  $e_i$  are random variables. It is assumed that the expected value of the error is zero, and that the error is uncorrelated with the two factors and with the errors of other assets. However, it is not assumed that the two factors are uncorrelated with each other. These factors are presumably observable variables, and their statistical properties can be studied independently of the asset returns.

In the case of the two-factor model we easily derive the following values for the expected rates of return and the covariances:

$$\bar{r}_i = a_i + b_{1i}\bar{f}_1 + b_{2i}\bar{f}_2$$

$$\text{cov}(r_i, r_i) = \begin{cases} b_{1i}b_{1j}\sigma_{f_1}^2 + (b_{1i}b_{2j} + b_{2i}b_{1j})\text{cov}(f_1, f_2) + b_{2i}b_{2j}\sigma_{f_2}^2, & i \neq j \\ b_{1i}^2\sigma_{f_1}^2 + 2b_{1i}b_{2i}\text{cov}(f_1, f_2) + b_{2i}^2\sigma_{f_2}^2 + \sigma_{e_i}^2, & i = j. \end{cases}$$

The  $b_{1i}$ 's and  $b_{2i}$ 's can be obtained by forming the covariance of  $r_i$  with  $f_1$  and  $f_2$ , leading to

$$\text{cov}(r_i, f_1) = b_{1i}\sigma_{f_1}^2 + b_{2i}\sigma_{f_1, f_2}$$

$$\text{cov}(r_i, f_2) = b_{1i}\sigma_{f_1, f_2} + b_{2i}\sigma_{f_2}^2.$$

These give two equations that can be solved for the two unknowns  $b_{1i}$  and  $b_{2i}$ .

A two-factor model is often an improvement of a single-factor model. For example, suppose a single-factor model were proposed and the  $a_i$ 's and  $b_i$ 's determined by fitting data. It might be found that the resulting error terms are large and that they exhibit correlation with the factor and with each other. In this case the single-factor model is not a good representation of the actual returns structure. A two-factor model may lead to smaller error terms, and these terms may exhibit the assumed correlation properties. The two-factor model will still be much simpler than a full unstructured covariance matrix.

It should be clear how to extend the model to include a greater number of factors. Quite comprehensive models of this type have been constructed. It is generally agreed that for models of U.S. stocks, it is appropriate to use between 3 and 15 factors.

## Selection of Factors

The selection of appropriate factors for a factor model is part science and part art (like most practical analyses). It is helpful, however, to place factors in three categories.

Once these categories are recognized, you will no doubt be able to dream up additional useful factors. Here are the categories:

**1. External factors** Very commonly, factors are chosen to be variables that are external to the securities being explicitly considered in the model. Examples are gross national product (GNP), consumer price index (CPI), unemployment rate, or a new construction index. The U.S. Government publishes numerous such statistics. It is possible to use other external variables as well, such as the number of traffic accidents in a month or sun spot activity.

**2. Extracted factors** It is possible to extract factors from the known information about security returns. For example, the factor used most frequently is the rate of return on the market portfolio. This factor is constructed directly from the returns of the individual securities. As another example, the rate of return of one security can be used as a factor for others. More commonly, an average of the returns of the securities in an industry is used as a factor; for example, there might be an industrial factor, a utilities factor, and a transportation factor. Factors can also be extracted by the method of principal components. (See Exercise 3.) This method uses the covariance matrix of the returns to find combinations of securities that have large variances. Indeed, extracted factors are usually linear combinations of the returns of individual securities (as in the preceding examples). Factors can be extracted in more complex ways. For example, a factor might be defined as the ratio of the returns of two stocks, the number of days since the last market peak, or a moving average of the market return.

**3. Firm characteristics** Firms are characterized financially by a number of firm-specific values, such as the price–earnings ratio, the dividend-payout ratio, an earnings forecast, and many other variables. About 50 such variables for each major security are available from various data services. These characteristics can be used in a factor model. The characteristics do not serve as factors in the usual sense, but they play a similar role. As an example, suppose that we decide to use a single factor  $f$  (of the normal kind) and a single firm characteristic  $g$  (such as last quarter's price–earnings ratio). We then represent the rate of return on security  $i$  as

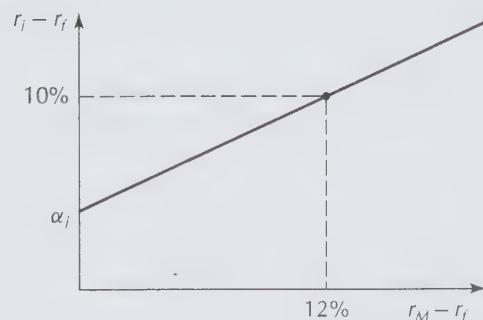
$$r_i = a_i + b_i f + c g_i + e_i. \quad (8.4)$$

In this model, the constant  $a$  is the same for each security, but  $g_i$  (the value of the characteristic) varies. The characteristic term does not contribute to systematic (or nondiversifiable) risk, but rather it may reduce the variance of the error term  $e_i$ . In other words, the term  $c g_i$  can be regarded as an estimate of the error term that would appear in the standard single-factor model. Firm characteristics are effective additions to factor models.

## 8.3 The CAPM as a Factor Model

The CAPM can be derived as a special case of a single-factor model. This view adds considerable insight to the CAPM development.

**FIGURE 8.2 Characteristic line.**  
This line represents a single-factor model that has  $r_M - r_f$  as the factor for the variable  $r_i - r_f$ .



## The Characteristic Line

Let us hypothesize a single-factor model for stock returns, with the factor being the market rate of return  $r_M$ . For convenience we can subtract the constant  $r_f$  from this factor and also from the rate of return  $r_i$ , thereby expressing the model in terms of the excess returns  $r_i - r_f$  and  $r_M - r_f$ . The factor model then becomes

$$r_i - r_f = \alpha_i + \beta_i(r_M - r_f) + e_i. \quad (8.5)$$

It is conventional to use the notation  $\alpha_i$  and  $\beta_i$  for the coefficients of this special model, rather than the  $a_i$ 's and  $b_i$ 's that are being used more generally. Again it is assumed that  $E(e_i) = 0$  and that  $e_i$  is uncorrelated with the market return (the factor) and with other  $e_j$ 's.

The **characteristic equation** or **characteristic line** corresponding to (8.5) is the line formed by putting  $e_i = 0$ ; that is, it is the line  $r_i - r_f = \alpha_i + \beta_i(r_M - r_f)$  drawn on a diagram of  $r_i$  versus  $r_M$ . Such a line is shown in Figure 8.2. A single typical point is indicated on the line. If measurements of  $r_i - r_f$  and  $r_M - r_f$  were taken and plotted on this diagram, they would fall at various places, but the characteristic line would presumably define a good fit through the scatter of points.

The expected value of this equation is

$$\bar{r}_i - r_f = \alpha_i + \beta_i(\bar{r}_M - r_f),$$

which is identical to the CAPM except for the presence of  $\alpha_i$ . The CAPM predicts that  $\alpha_i = 0$ .

The value of  $\beta_i$  in this model can be calculated directly. We take the covariance of both sides of (8.5) with  $r_M$ . This produces

$$\sigma_{iM} = \beta_i \sigma_M^2$$

and hence

$$\beta_i = \frac{\sigma_{iM}}{\sigma_M^2}.$$

This is exactly the same expression that holds for the  $\beta_i$  used in the CAPM (and that is why we use the same notation).

The characteristic line is in a sense more general than the CAPM because it allows  $\alpha_i$  to be nonzero. From the CAPM viewpoint,  $\alpha_i$  can be regarded as a measure of the amount that asset  $i$  is mispriced. A stock with a positive  $\alpha_i$  is, according to this view, performing better than it should, and a stock with a negative  $\alpha_i$  is performing worse than it should. Some financial services organizations (and some highly technical investors) estimate  $\alpha$  as well as  $\beta$  for a large assortment of stocks. Note, however, that the single-factor model that leads to the CAPM formula is not equivalent to the general model underlying the CAPM, since the general model is based on an arbitrary covariance matrix, but assumes that the market is efficient. The single-factor model has a very simple covariance structure, but makes no assumption about efficiency.

**Example 8.2 (Four stocks and the market)** Let us rework Example 8.1 by using the excess market return as a factor. We assume that the market consists of just the four stocks, with equal weights. Therefore the market return in any year is just the average of the returns of the four stocks. These are shown in the upper portion of Table 8.2. We also adjoin the historical value of the risk-free rate of return for each of the 10 years. The relevant statistical quantities are computed by the same estimating formulas as in the earlier example, except that the factor is taken to be the excess return on the market, which will change the formula for  $a_i$  to  $\alpha_i$ . As seen from the table, a large portion of the variability of the stock returns is explained by the factor. In other words, there is relatively low nonsystematic risk. Furthermore, a side calculation shows that the errors are close to being uncorrelated with each other and

**TABLE 8.2**  
**FACTOR MODEL WITH MARKET**

Year	Stock 1	Stock 2	Stock 3	Stock 4	Market	Riskless
1	11.91	29.59	23.27	27.24	23.00	6.20
2	18.37	15.25	19.47	17.05	17.54	6.70
3	3.64	3.53	-6.58	10.20	2.70	6.40
4	24.37	17.67	15.08	20.26	19.34	5.70
5	30.42	12.74	16.24	19.84	19.81	5.90
6	-1.45	-2.56	-15.05	1.51	-4.39	5.20
7	20.11	25.46	17.80	12.24	18.90	4.90
8	9.28	6.92	18.82	16.12	12.78	5.50
9	17.63	9.73	3.05	22.93	13.34	6.10
10	15.71	25.09	16.94	3.49	15.31	5.80
aver	15.00	14.34	10.90	15.09	13.83	5.84
var	90.28	107.24	162.19	68.27	72.12	
cov	65.08	73.62	100.78	48.99	72.12	
$\beta$	.90	1.02	1.40	.68	1.00	
$\alpha$	1.95	.34	-6.11	3.82	0.00	
e-var	31.54	32.09	21.37	34.99		

Now the factor is taken to be the excess return on the market portfolio. The variation in stock returns is largely explained by this return, and the errors are uncorrelated with each other and with the market. This model provides an excellent fit to the data.

with the market return. For example, the data provide the estimates  $\text{cov}(e_1, e_2) = -14$  and  $\text{cov}(e_2, e_3) = 2$ , which are much smaller than for the earlier model. We conclude that this single-factor model is an excellent representation of the stock returns of the four stocks. In other words, for this example, the market return serves as a much better factor than the industrial index factor used earlier. However, this may not be true for other examples.

## 8.4 Arbitrage Pricing Theory\*

The factor model framework leads to an alternative theory of asset pricing, termed **arbitrage pricing theory** (APT). This theory does not require the assumption that investors evaluate portfolios on the basis of means and variances; only that, when returns are certain, investors prefer greater return to lesser return. In this sense the theory is much more satisfying than the CAPM theory, which relies on both the mean–variance framework and a strong version of equilibrium, which assumes that everyone uses the mean–variance framework.

The APT does, however, require a special assumption of its own. This is the assumption that the universe of assets being considered is large. For the theory to work exactly, we must, in fact, assume that there are an infinite number of securities, and that these securities differ from each other in nontrivial ways. This assumption is generally felt to be satisfied well enough by, say, the universe of all publicly traded U.S. stocks.

### Simple Version of APT

To explain the concept underlying the APT, we first consider an idealized special case. Assume that all asset rates of return satisfy the following one-factor model:

$$r_i = a_i + b_i f.$$

Different assets will have different  $a_i$ 's and  $b_i$ 's. This factor model is special because there is no error term. The uncertainty associated with a return is due only to the uncertainty in the factor  $f$ . The point of APT is that the values of  $a_i$  and  $b_i$  must be related if arbitrage opportunities are to be excluded. To work out the relationship between  $a_i$  and  $b_i$  we write the model for two assets  $i$  (as before) and  $j$ , which is

$$r_j = a_j + b_j f.$$

The only requirement in the selection of these two securities is that  $b_i \neq b_j$ . Now form a portfolio with weights  $w_i = w$  and  $w_j = 1 - w$ . We know that the rate of return of this portfolio is

$$r = wa_i + (1-w)a_j + [wb_i + (1-w)b_j]f.$$

We shall select  $w$  so that the coefficient of  $f$  in this equation is zero. Specifically, we select  $w = b_j/(b_j - b_i)$ . This yields a rate of return of

$$r = wa_i + (1 - w)a_j = \frac{a_i b_j}{b_j - b_i} + \frac{a_j b_i}{b_i - b_j}. \quad (8.6)$$

This special portfolio is risk free because the equation for  $r$  contains no random element. If there is a separate risk-free asset with rate of return  $r_f$ , it is clear that the portfolio constructed in (8.6) must have this same rate—otherwise there would be an arbitrage opportunity. Even if there is no explicit risk-free asset, all portfolios constructed this way, with no dependence on  $f$ , must have the same rate of return. We denote this rate by  $\lambda_0$ , recognizing that  $\lambda_0 = r_f$  if there is an explicit risk-free asset.

Setting the right-hand side of (8.6) equal to  $\lambda_0$ , we find

$$\lambda_0(b_j - b_i) = a_i b_j - a_j b_i,$$

which can be rearranged to

$$\frac{a_j - \lambda_0}{b_j} = \frac{a_i - \lambda_0}{b_i}.$$

This is a general relation that must hold for all  $i$  and  $j$ . Therefore,

$$\frac{a_i - \lambda_0}{b_i} = c$$

holds for all  $i$  for some constant  $c$ . This shows explicitly that the values of  $a_i$  and  $b_i$  are not independent. Indeed,  $a_i = \lambda_0 + b_i c$ .

To see that such a relation is reasonable, suppose we take  $f$  to be the rate of return on the S&P 500 average. If  $a_i$  and  $b_i$  were arbitrary, we might specify a stock  $i$  with  $a_i = .50$  and  $b_i = 1.0$ , which would give  $i$  a rate of return of 50% plus the S&P 500 rate. Clearly this is unreasonably high. No stock does this well. More realistically, if we have  $a_i = .50$ , then  $b_i$  will be negative so that, overall,  $r_i$  makes sense. As another case, if  $a_i$  is the risk-free rate, then  $b_i$  should be zero. The relation  $a_i = \lambda_0 + b_i c$  keeps things in proper alignment.

We can use this information to write a simple formula for the expected rate of return of asset  $i$ . We have

$$\bar{r}_i = a_i + b_i \bar{f} = \lambda_0 + b_i c + b_i \bar{f}$$

or, alternatively,

$$\bar{r}_i = \lambda_0 + b_i \lambda_1 \quad (8.7)$$

for the constant  $\lambda_1 = c + \bar{f}$ . We see that once the constants  $\lambda_0$  and  $\lambda_1$  are known, the expected return of an asset is determined entirely by the factor loading  $b_i$  (since  $a_i$  must follow  $b_i$ ).

Notice that pricing formula (8.7) looks similar to the CAPM. If the factor  $f$  is chosen to be the rate of return on the market  $r_M$ , then we can set  $\lambda_0 = r_f$  and  $\lambda_1 = \bar{r}_M - r_f$ , and the APT is identical to the CAPM with  $b_i = \beta_i$ .

For additional factors the result is similar. We now give a more general statement and proof:

**Simple APT** Suppose that there are  $n$  assets whose rates of return are governed by  $m < n$  factors according to the equation

$$r_i = a_i + \sum_{j=1}^m b_{ij} f_j$$

for  $i = 1, 2, \dots, n$ . Then there are constants  $\lambda_0, \lambda_1, \dots, \lambda_m$  such that

$$\bar{r}_i = \lambda_0 + \sum_{j=1}^m b_{ij} \lambda_j$$

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

**Proof:** We prove the statement for the case of two factors. Suppose we invest a dollar amount  $x_i$  in asset  $i$ ,  $i = 1, 2, \dots, n$ , in order to satisfy  $\sum_{i=1}^n x_i = 0$ ,  $\sum_{i=1}^n x_i b_{i1} = 0$ , and  $\sum_{i=1}^n x_i b_{i2} = 0$ . This portfolio requires zero net investment and has zero risk. Therefore its expected payoff must be zero. Hence  $\sum_{i=1}^n x_i \bar{r}_i = 0$ . Defining the vectors  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ ,  $\mathbf{b}_1 = (b_{11}, b_{21}, \dots, b_{n1})$ ,  $\mathbf{b}_2 = (b_{12}, b_{22}, b_{32}, \dots, b_{n2})$ ,  $\mathbf{1} = (1, 1, \dots, 1)$ , and  $\bar{\mathbf{r}} = (\bar{r}_1, \bar{r}_2, \dots, \bar{r}_n)$ , we can restate the foregoing as follows: For any  $\mathbf{x}$  satisfying  $\mathbf{x}^T \mathbf{1} = 0$ ,  $\mathbf{x}^T \mathbf{b}_1 = 0$ , and  $\mathbf{x}^T \mathbf{b}_2 = 0$  it follows that  $\mathbf{x}^T \bar{\mathbf{r}} = 0$ ; that is, any  $\mathbf{x}$  orthogonal to  $\mathbf{1}$ ,  $\mathbf{b}_1$ , and  $\mathbf{b}_2$  is also orthogonal to  $\bar{\mathbf{r}}$ . It follows from a standard result in linear algebra<sup>2</sup> that  $\bar{\mathbf{r}}$  must be a linear combination of the vectors  $\mathbf{1}$ ,  $\mathbf{b}_1$ , and  $\mathbf{b}_2$ . Thus there are constants  $\lambda_0, \lambda_1, \lambda_2$  such that  $\bar{\mathbf{r}} = \lambda_0 \mathbf{1} + \mathbf{b}_1 \lambda_1 + \mathbf{b}_2 \lambda_2$ . This is identical to the given statement. ■

To understand this result, let us look at some special cases. If all the  $b_{ij}$ 's are zero, then there is no risk and we have  $a_i = \lambda_0$ , which is appropriate. If a  $b_{ij}$  is nonzero, then  $\bar{r}_i$  increases in proportion to  $b_{ij}$ ; the value  $\lambda_j$  is the **price of risk** associated with the factor  $f_i$ , often called the **factor price**. As one accepts greater amounts of  $f_i$ , one obtains greater expected return.

## Well-Diversified Portfolios

We now consider more realistic factor models, which have error terms as well as factor terms. Suppose there are a total of  $n$  assets and the rate of return on asset  $i$  satisfies

$$r_i = a_i + \sum_{j=1}^m b_{ij} f_j + e_i$$

---

<sup>2</sup> You can visualize this in the three dimensions of a room. Fix a vector  $\mathbf{b}$ , say, running along the floor and perpendicular to a wall. Suppose that for all  $\mathbf{x}$  with  $\mathbf{x}^T \mathbf{b} = 0$ , there also holds  $\mathbf{x}^T \bar{\mathbf{r}} = 0$ . The set of  $\mathbf{x}$ 's are those on the wall. Then you should see that  $\bar{\mathbf{r}} = \lambda \mathbf{b}$  for some  $\lambda$ .

where  $E(e_i) = 0$  and  $E[e_i]^2 = \sigma_{e_i}^2$ . Also assume that  $e_i$  is uncorrelated with the factors and with the error terms of other assets. Let us form a portfolio using the weights  $w_1, w_2, \dots, w_n$  with  $\sum_{i=1}^n w_i = 1$ . The rate of return of the portfolio is

$$r = a + \sum_{j=1}^m b_j f_j + e$$

where

$$a = \sum_{i=1}^n w_i a_i$$

$$b_j = \sum_{i=1}^n w_i b_{ij}$$

$$\sigma_e^2 = \sum_{i=1}^n w_i^2 \sigma_{e_i}^2.$$

Suppose that for each  $i$  there holds  $\sigma_{e_i}^2 \leq S^2$  for some constant  $S$ . Suppose also that the portfolio is **well diversified** in the sense that for each  $i$  there holds  $w_i \leq W/n$  for some constant  $W \approx 1$ . This assures that no one asset is heavily weighted in the portfolio. We then find that

$$\sigma_e^2 \leq \frac{1}{n^2} \sum_{i=1}^n W^2 S^2 \leq \frac{1}{n} W^2 S^2.$$

We now let  $n \rightarrow \infty$ . While doing this we assume that the bound  $\sigma_{e_i}^2 \leq S^2$  remains valid for all  $i$ . Also for each  $n$ , we select a portfolio that is well diversified. As  $n \rightarrow \infty$ , we see that  $\sigma_e^2 \rightarrow 0$ . In other words, the error term associated with a well-diversified portfolio of an infinite number of assets has a variance of zero. For a finite, but large, number of assets the error term has approximately zero variance.

## General APT

We now combine the ideas of the preceding two subsections. We imagine forming thousands of different well-diversified portfolios, each being (essentially) error free. These portfolios form a collection of assets, the return on each satisfying a factor model without error. We therefore can apply the simple APT to conclude that there are constants  $\lambda_0, \lambda_1, \dots, \lambda_m$  such that for any well-diversified portfolio having a rate of return

$$r = a + \sum_{j=1}^m b_j f_j$$

the expected rate of return is

$$\bar{r} = \lambda_0 + \sum_{j=1}^m b_j \lambda_j.$$

Since various well-diversified portfolios can be formed with weights that differ on only a small number of basic assets, it follows that these individual assets must also satisfy

$$\bar{r}_i = \lambda_0 + \sum_{j=1}^m b_{ij}\lambda_j.$$

(This argument is not completely rigorous; but a more rigorous argument is quite complex.)

This is again basically a relation that says that  $a_i$  is not independent of the  $b_{ij}$ 's. The risk-free term must be related to the factor loadings. This is true even when there are error terms, provided there is a large number of assets so that error terms can be effectively diversified away.

## APT and CAPM

The factor model underlying APT can be applied to the CAPM framework to derive a relation between the two theories.

Using a two-factor model we have

$$r_i = a_i + b_{i1}f_1 + b_{i2}f_2 + e_i.$$

We find the covariance of this asset with the market portfolio to be

$$\text{cov}(r_M, r_i) = b_{i1}\text{cov}(r_M, f_1) + b_{i2}\text{cov}(r_M, f_2) + \text{cov}(r_M, e_i).$$

If the market represents a well-diversified portfolio, it will contain essentially no error term, and hence it is reasonable to ignore the term  $\text{cov}(r_M, e_i)$  in the foregoing expression. We can then write the beta of the asset as

$$\beta_i = b_{i1}\beta_{f_1} + b_{i2}\beta_{f_2},$$

where

$$\begin{aligned}\beta_{f_1} &= \sigma_{M,f_1}/\sigma_M^2 \\ \beta_{f_2} &= \sigma_{M,f_2}/\sigma_M^2.\end{aligned}$$

Hence the overall beta of the asset can be considered to be made up from underlying factor betas that do not depend on the particular asset. The weight of these factor betas in the overall asset beta is equal to the factor loadings. Hence in this framework, the reason that different assets have different betas is that they have different loadings.

## 8.5 Projection Pricing with Factors

The factor framework can be useful when one seeks to price an asset directly, rather than compute an implied expected rate of return. The linear structure of factors is a natural setting when assigning prices since prices are combined linearly in a market, and assignment of price is often a primary objective when analyzing a new asset.

Suppose there are  $n$  such assets that span a significant submarket described by factors. In a manner similar to a standard factor framework, each asset  $i$  in that submarket has payoff

$$y_i = a_i + \sum_{j=1}^m b_{ij} f_j + \varepsilon_i,$$

where  $y_i$  is the payoff of the  $i$ th asset. The factors  $f_j$  are random and represent the uncertain influence of outside (or inside) random elements. For instance, a factor may be the cost of energy. The  $b_{ij}$ 's are the factor loadings for  $i$  of the factors  $j$ . The  $\varepsilon_i$ 's are the individual errors and are assumed to be uncorrelated with each other and with the factors.

A new asset has payoff  $y$  given by

$$y = a + \sum_{j=1}^m b_j f_j + \varepsilon,$$

where, consistent with the structure of the  $n$  given assets, we assume that  $\varepsilon$  is independent of all other uncertainties.

To price this asset by projection we use the (scaled) minimum-norm portfolio of marketed assets having unit price. This scaled portfolio, being a linear combination of the original assets, is of the form

$$y^q = a^q + \sum_{j=1}^m b_j^q f_j + \varepsilon^q.$$

The projection price of  $y$  is then found by taking the expectation of the product, namely,  $P_y = E[y^q y]$ . This works out to be

$$P_y = \theta_0 a + \sum_{j=1}^m \theta_j b_j, \quad (8.8)$$

where

$$\begin{aligned} \theta_0 &= a^q + \sum_{j=1}^m b_j^q \bar{f}_j \\ \theta_j &= a^q \bar{f}_j + \sum_{i=1}^m b_i^q \bar{f}_i \bar{f}_j. \end{aligned} \quad (8.9)$$

We see from equation (8.8) that each  $\theta_j$  can be thought of as a unit price for its corresponding factor. For example, for an energy factor, the corresponding  $\theta$  is a value of a unit of the energy factor loading.<sup>3</sup>

---

<sup>3</sup> If the error  $e$  of the new asset is not uncorrelated with the errors of the other assets, there will be additional terms. These go to zero if the minimum-norm portfolio is well diversified.

**Example 8.3 (An exercise to try)** Suppose that a submarket is characterized by a single factor  $f$  such that a typical asset has payoff of the form

$$y = c + bf + \varepsilon. \quad (8.10)$$

The values of the constants  $c$  and  $b$  are associated with a particular asset, and  $\varepsilon$  is a random error. The expected values of  $f$  and  $\varepsilon$  are zero, and their variances are  $\sigma_f^2$  and  $\sigma_\varepsilon^2$ . Further,  $\varepsilon$  is uncorrelated with  $f$  and with the errors of all other assets in the market.

Suppose that the market is spanned by just two assets. Each has price 1, and the two corresponding payoffs are

$$y_1 = 1 \cdot f + \varepsilon_1 \quad (8.11)$$

$$y_2 = R. \quad (8.12)$$

Here  $\varepsilon_1$  is a random error. The factor loading is 1. The expected values of  $f$  and  $\varepsilon_1$  are zero, and their variances are  $\sigma_f^2$  and  $\sigma_{\varepsilon_1}^2$ . Further,  $\varepsilon_1$  is uncorrelated with  $f$ .

- (a) Find the weights  $w_1 = w$  and  $w_2 = 1 - w$  for the minimum-norm portfolio  $y_{\min}$  of these two assets. [You should find  $w = R^2 / (\sigma_f^2 + \sigma_{\varepsilon_1}^2 + R^2)$ .]
- (b) Scale this portfolio by a multiple  $s$  so that it prices  $y_2$  correctly, that is,  $E[sy_{\min}R] = 1$ . [You should get  $s = 1 / ((1 - w)R^2)$ .]
- (c) Evaluate the general asset (8.10) as  $P_y = E[sy_{\min}y]$  to obtain the factor pricing equation

$$P_y = \theta_0 c + \theta_1 b.$$

[You should obtain  $\theta_0 = 1/R$  and  $\theta_1 = \sigma_f^2 / (\sigma_f^2 + \sigma_\varepsilon^2)$ .]

## 8.6 A Multiperiod Fallacy

The CAPM theory is a beautiful and simple theory that follows very logically from the single-period mean-variance theory of Markowitz. In practice, however, both mean-variance theory and the derived CAPM are applied to situations that are inherently multiperiod, such as the construction of portfolios of common stocks that can be traded at any time.

The simplest way to apply mean-variance theory to the multiperiod case is to select a basic period length—say, 1 month. The Markowitz problem is formulated for this period. If this problem is solved, it should, according to the CAPM assumption, prescribe that the optimal portfolio weighting vector  $\mathbf{w}$  is equal to the market portfolio. This idea can then be carried forward another period. If it is assumed that the statistical properties of the returns for the next period are identical to those of the previous period and the new returns are uncorrelated with those of the previous period, the new weighting vector  $\mathbf{w}$  will be equal to that of the previous period. However, in the meantime the prices will have changed relative to each other; and hence the vector

$w$  will no longer correspond to the market portfolio since the market weights are capitalization weights, and a price variation changes the capitalization. This is a basic fallacy, or contradiction, since the Markowitz model keeps giving the same weights, but the market portfolio weights change every period.

Let us consider a simple example. Suppose that there are only two stocks, each having the same initial price of, say, \$1, the same mean and variance of return, and zero correlation with each other. Both stocks are in equal supply in the market—say, 1,000 shares of each. Suppose that we have an amount  $X_0$  to invest. By symmetry, the mean-variance solution will be  $w = (\frac{1}{2}, \frac{1}{2})$ ; hence we should purchase equal amounts of both assets (equal dollar amounts, which is equivalent to equal numbers of shares since the prices of the two stocks are equal). This solution corresponds to the market portfolio.

Suppose that during the first period the first stock doubles in value and the second does not change. Hence now  $p_1 = \$2$  and  $p_2 = \$1$ , and our total wealth has increased to  $1.5X_0$ . Since the statistical properties remain unchanged, the optimal mean-variance solution will still have  $w = (\frac{1}{2}, \frac{1}{2})$ . This implies that we should again divide our money evenly between the two stocks. But if we do that we will purchase  $\frac{1}{4}1.5X_0$  shares of stock 1 and  $\frac{1}{2}1.5X_0$  shares of stock 2. This does *not* correspond to the market portfolio, which still has equal numbers of shares of the two stocks. In general, as prices change relative to each other, the dollar proportions represented in the market also change; but a repeating mean-variance model dictates that the dollar proportions of an optimal portfolio should remain fixed, which is a contradiction.

The fallacy can be repaired by assuming that the expected returns change each period in a way that keeps the market portfolio optimal; but this destroys the elegance of the model. It is more satisfying to develop a full multiperiod approach (as in Part 4 of this text). The multiperiod approach reverses some conclusions of the single-period theory. For example, the multiperiod theory suggests that price volatility is actually desirable, rather than undesirable. Nevertheless, the single-period framework of Markowitz and the CAPM are beautiful theories that ushered in an era of quantitative analysis and have provided an elegant foundation to support further work.

## 8.7 Summary

Special analytical procedures and modeling techniques can make mean-variance portfolio theory more practical than it would be if the theory were used in its barest form. The procedures and techniques discussed in this chapter include: (1) factor models to reduce the number of parameters required to specify a mean-variance structure, (2) use of APT to add factors to the CAPM and also to avoid the equilibrium assumption that underlies the CAPM, and (3) recognition of the errors inherent in computing parameter estimates from historical records of returns.

A factor model expresses the rate of return of each asset as a linear combination of certain specified (random) factor variables. The same factors are used for each asset, but the coefficients of the linear combination of these factors are different for different assets. In addition to the factor terms, there are a constant term  $a_i$  and an error term

$e_i$ . The coefficients of the factors are called factor loadings. In making calculations with the model, it is usually assumed that the error terms are uncorrelated with each other and with the factors.

A great advantage of a factor model is that it has far fewer parameters than a standard mean–variance representation. In practice, between three and fifteen factors can provide a good representation of the covariance properties of the returns of thousands of U.S. stocks.

There are several choices for factors. The most common choice is the return on the market portfolio. A factor model using this single factor is closely related to the CAPM. Other choices include various economic indicators published by the U.S. government or factors extracted as combinations of certain asset returns. It is also helpful to supplement a factor model by including combinations of company-specific financial characteristics.

When the excess market return is used as the single factor, the resulting factor model can be interpreted as defining a straight line on a graph with  $r_M - r_f$  being the horizontal axis and  $r - r_f$  the vertical axis. This line is called the characteristic line of the asset. Its vertical intercept is called alpha, and its slope is the beta of the CAPM. The CAPM predicts that alpha is zero (but in practice it may be nonzero).

Arbitrage pricing theory (APT) is built directly on a factor model. For the theory to be useful, it is important that the underlying factor model be a good representation in the sense that the error terms are uncorrelated with each other and with the factors. In that case, the error terms can be diversified away by forming combinations of a large number of assets.

The result of APT is that the coefficients of the underlying factor model must satisfy a linear relation. In the special case where the underlying factor model has the single factor equal to the excess return on the market portfolio, the CAPM theory states that  $\alpha = 0$ . This is a special case of APT, which states that the constant  $a$  in the expression for the return of an asset is a linear combination of the factor loadings of that asset. Again, the difficult part of applying APT is the determination of appropriate factors.

It is tempting to assume that the parameter values necessary to implement mean–variance theory—the expected returns, variances, and covariances for a Markowitz formulation, or the  $a_i$ 's and  $b_{ij}$ 's for a factor model representation—can be estimated from historical returns data. As discussed in the next chapter, some parameter values can be estimated this way, but others cannot.

The Markowitz mean–variance formulation of portfolio theory and the subsequent theories of CAPM, factor models, and APT provide an elegant foundation for single-period investment analysis. These developments have elaborated the benefits of diversification and deepened our understanding of risk in a market environment. These theories have also provided approaches that can be implemented. Indeed, this whole area has had a profound influence on the practice of portfolio management: index funds now abound, betas are computed and widely discussed in the financial community, large quadratic programming programs have been written to solve the Markowitz problem, numerous factor models have been constructed and tested, and trillions of dollars have been managed with at least some guidance from these ideas and methods.

But mean–variance theory is not a universal investment panacea. The assumption that all investors focus exclusively on mean and variance is questionable; it is hard to estimate the required parameter values, it seems unlikely (as required of the equilibrium argument) that everyone has the same estimates of the parameter values, and the approach must be modified in a multiperiod framework. Each of these difficulties can be overcome to some extent by extending the model, living with approximations, or looking deeper into the properties of the assets under consideration. A great deal of innovative effort has been so devoted. But ultimately, to make significant progress, we must expand the fundamental tools of analysis beyond mean–variance. We must formulate a theory that, built on the insights of the mean–variance approach, treats uncertainty more explicitly and is directed at multiperiod situations.

## Exercises

- 1. (A simple portfolio)** Someone who believes that the collection of all stocks satisfies a single-factor model with the market portfolio serving as the factor gives you information on three stocks which make up a portfolio. (See Table 8.3.) In addition, you know that the market portfolio has an expected rate of return of 12% and a standard deviation of 18%. The risk-free rate is 5%.
- What is the portfolio's expected rate of return?
  - Assuming the factor model is accurate, what is the standard deviation of this rate of return?

**TABLE 8.3  
SIMPLE PORTFOLIO**

Stock	Beta	Standard deviation of random error term	Weight in portfolio
A	1.10	7.0%	20%
B	0.80	2.3%	50%
C	1.00	1.0%	30%

- 2. (APT factors)** Two stocks are believed to satisfy the two-factor model

$$r_1 = a_1 + 2f_1 + f_2$$

$$r_2 = a_2 + 3f_1 + 4f_2.$$

In addition, there is a risk-free asset with a rate of return of 10%. It is known that  $\bar{r}_1 = 15\%$  and  $\bar{r}_2 = 20\%$ . What are the values of  $\lambda_0$ ,  $\lambda_1$ , and  $\lambda_2$  for this model?

- 3. (Principal components ⊕)** Suppose there are  $n$  random variables  $x_1, x_2, \dots, x_n$  and let  $\mathbf{V}$  be the corresponding covariance matrix. An **eigenvector** of  $\mathbf{V}$  is a vector  $\mathbf{v} = (v_1, v_2, \dots, v_n)$  such that  $\mathbf{V}\mathbf{v} = \lambda\mathbf{v}$  for some  $\lambda$  (called an eigenvalue of  $\mathbf{V}$ ). The random variable  $v_1x_1 + v_2x_2 + \dots + v_nx_n$  is a **principal component**. The first principal component is the one corresponding to the largest eigenvalue of  $\mathbf{V}$ , the second to the second largest, and so forth.

A good candidate for the factor in a one-factor model of  $n$  asset returns is the first principal component extracted from the  $n$  returns themselves; that is, by using the principal eigenvector of the covariance matrix of the returns. Find the first principal component for

the data of Example 8.2. Does this factor (when normalized) resemble the return on the market portfolio? [Note: For this part, you need an eigenvector calculator as available in most matrix operations packages.]

4. (Variance estimate) Let  $r_i$ , for  $i = 1, 2, \dots, n$ , be independent samples of a return  $r$  of mean  $\bar{r}$  and variance  $\sigma^2$ . Define the estimates

$$\hat{r} = \frac{1}{n} \sum_{i=1}^n r_i$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (r_i - \hat{r})^2.$$

Show that  $E(s^2) = \sigma^2$ .

5. (Find loading) The data here show the rate of return of a stock and the corresponding value of a factor  $F$ . In the model

$$r_i = a + b \times F + \varepsilon_i$$

find the best values of  $a$  and  $b$  (which minimize the total squared error). A simple way is to optimize with a spreadsheet program.

Rate	.02	.34	.12	.34	.67	.11	.33	.56	.34	-.28
Factor F	.40	.80	.14	.66	.90	.20	.18	.47	.21	-.30

6. (3 portfolio APT) Assume that the following two-index model describes returns:

$$\bar{r}_i = r_f + b_{i1}\lambda_1 + b_{i2}\lambda_2.$$

Assume that the following three portfolios are observed:

Portfolio	Expected Return(%)	$b_{i1}$	$b_{i2}$
A	10	1	0
B	1	0	2
C	-2	-1	1

- (a) According to APT, what is the expected return, factor-loading relationship for this market.

- (b) Now consider portfolio D with the following characteristics:

$$r_D = 15\%, \quad b_{D1} = 2, \quad b_{D2} = 1.$$

Is an arbitrage opportunity possible? If yes, then describe the arbitrage opportunity.

7. (Cancellation) Suppose that sixteen stocks have been identified whose rates of return satisfy

$$\bar{r}_i = \pm\alpha + f + \varepsilon_i,$$

where  $\alpha > 0$ . Eight of the stocks use the  $+$  sign and the other eight use the  $-$  sign. The factor  $f$  is common to all sixteen stocks. It has mean equal to 1, and its standard deviation

is 15%. Each  $\varepsilon_i$  represents firm specific error, in the sense that each has zero mean, zero covariance with  $f_j$ , and zero covariance with other stocks. Each  $\varepsilon_i$  has a standard deviation of 24%. Now assume that a portfolio consists of all of these stocks, with equal weight given to each one of them. What is the expected rate of return and the corresponding standard deviation of that rate?

## References

The factor analysis approach to structuring a family of returns is quite well developed. A good survey is contained in [1]. Also see [2]. The APT was devised by Ross [3]. For a practical application see [4]. For introductory presentations of factor models and the APT consult the finance textbooks listed as references for Chapter 2.

1. Sharpe, W. F. (1982), "Factors in New York Stock Exchange Security Returns 1931–1979," *Journal of Portfolio Management*, **8**, Summer, 5–19.
2. King, B. F. (1966), "Market and Industry Factors in Stock Price Behavior," *Journal of Business*, **39**, January, 137–170.
3. Ross, S. A. (1976), "The Arbitrage Theory of Capital Asset Pricing," *Journal of Economic Theory*, **13**, 341–360.
4. Chen, N. F., R. Roll, and S. A. Ross (1986), "Economic Forces and the Stock Market," *Journal of Business*, **59**, 383–403.

# DATA AND STATISTICS

**A** major issue in the application of mean–variance portfolio theory to stocks is the estimation of the parameter values that the theory requires: the mean values of the returns of each of the stocks, the corresponding variances, and the covariances between them. These parameter values are not readily available, nor can they be deduced by logic as with physical processes such as a wheel of fortune, a coin flip, or the roll of die with clear payoffs and associated probabilities. This chapter discusses how mean–variance parameters can be estimated.

The estimates obtained are rarely perfect, and it is therefore important to quantify the impact that the use of imperfect estimates can have on the performance of a portfolio's performance. This knowledge may guide us to modify the design so that it is less susceptible to estimation errors. This chapter discusses this important issue.

## 9.1 Basic Estimation Methods

The basic method for estimating parameters for mean–variance portfolio design is to use historical data. It is a convenient method since suitable sources of data are readily available. Some financial service organizations either supply the data or provide the parameter estimates based on the data. The method is also reasonably reliable for certain of the parameters such as the variances and covariances; but it is decidedly *unreliable* for other parameters, such as the expected returns. The lack of reliability is not due to faulty data or difficult computation, it is due to a fundamental limitation

of the process of extracting estimates from data. It is a statistical limitation, which we loosely term the **blur of history**. It is important to understand the basic statistics of data processing and this fundamental limitation.

## Period-Length Effects

Suppose that the yearly return of a stock is  $1 + r_y$ . This yearly return can be considered to be the result of 12 monthly returns and thus can be written as

$$1 + r_y = (1 + r_1)(1 + r_2) \cdots (1 + r_{12}).$$

In this equation the monthly returns are *not* measured in yearly terms; they are the actual returns for the month. For small values of the  $r_i$ 's we can expand the product and keep only the first-order terms, as

$$1 + r_y \approx 1 + r_1 + r_2 + \cdots + r_{12}. \quad (9.1)$$

In other words,  $r_y \approx \sum_{i=1}^n r_i$ , which means that the yearly rate of return is approximately equal to the sum of the 12 individual monthly returns. This approximation ignores the compounding effect, but it is good enough for our present purpose, which is to estimate the rough magnitudes of the parameters.

Assume that the monthly returns of a given stock all have the same statistical properties and are mutually uncorrelated; that is, each monthly  $r_i$  has the same expected value  $\bar{r}$  and the same variance  $\sigma^2$ . Using approximation (9.1) we find that

$$\bar{r}_y = 12\bar{r}.$$

Likewise, we find

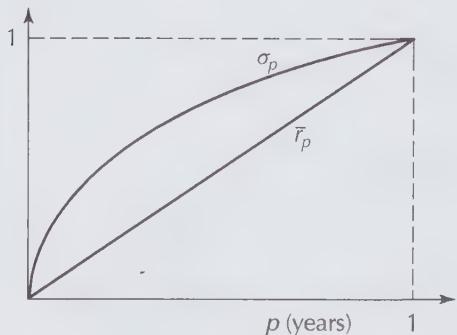
$$\sigma_y^2 = E \left[ \sum_{i=1}^{12} (r_i - \bar{r}) \right]^2 = E \left[ \sum_{i=1}^{12} (r_i - \bar{r})^2 \right] = 12\sigma^2,$$

where in the second step we used the fact that the returns are uncorrelated. Turning these equations around and taking the square root of the variance, we obtain an expression for the monthly values in terms of the yearly values,

$$\begin{aligned}\bar{r} &= \frac{1}{12}\bar{r}_y \\ \bar{\sigma} &= \frac{1}{\sqrt{12}}\sigma_y.\end{aligned}$$

This analysis can be generalized to any length of period, such as a week or a day. If we assume that the returns in different (identical length) periods have identical statistical properties and are uncorrelated, we obtain a similar result. Specifically, if the period is  $p$  part of a year (expressed as a fraction of a year), then the expected return and the standard deviation of the  $l$ -period rate of return can be found by generalizing from

**FIGURE 9.1 Period effects.** The expected rate of return over a period increases approximately linearly with the length of the period. The standard deviation increases as the square root of the length of the period.



monthly periods where  $p = 1/12$ . We have for general  $p$

$$\bar{r}_p = p\bar{r}_y \quad (9.2a)$$

$$\sigma_p = \sqrt{p}\sigma_y. \quad (9.2b)$$

It is the square-root term that causes the difficulty in estimation problems, as we shall see.

The effect of the period length on the expected rate of return and the standard deviation of the period returns is shown in Figure 9.1. The values for a 1-year period are normalized to unity for both the expected rate of return and the standard deviation. As the period is reduced, both the expected rate of return and the standard deviation of the period returns decrease. The expected rate of return is directly proportional to the length of the period. However, the standard deviation is proportional to the square root of the length of the period. This means that the ratio of the two—the ratio of standard deviation to expected rate of return—*increases* dramatically as the period length is reduced. In fact, this ratio goes to infinity as the period length goes to zero. Therefore the rates of return for small periods have very high standard deviations compared to their expected values.

Let us apply this analysis to a typical stock. The mean yearly rate of return for stocks ranges from around 6% to 30%, with a typical value being about 12%. These mean values change with time, so any particular value is meaningful only for about 2 or 3 years. The standard deviation of yearly stock returns ranges from around 10% to 60%, with 15% being somewhat typical.

Now let us translate the values of mean and variance into corresponding monthly values. Accordingly, we set  $p = 1/12$  in the formulas (9.2a) and (9.2b). Let us use the nominal values of  $\bar{r}_y = 12\%$  for the yearly expected rate of return, and  $\sigma_y = 15\%$  for the yearly standard deviation. This leads to  $\bar{r}_{1/12} = 1\%$  and  $\sigma_{1/12} = 4.33\%$  for the corresponding monthly values. Hence the standard deviation of the monthly return is 4.3 times the expected rate of return, whereas for the yearly figures the ratio is 1.25. The relative error is amplified as the period is shortened. Let us go a bit further and assume that returns are generated through independent *daily* returns. Assuming 250 trading days per year, we set  $p = 1/250$ . Then  $\bar{r}_{1/250} = .048\%$  and  $\sigma_{1/250} = .95\%$  are the corresponding daily values. The ratio of the two is now  $.95/.048 = 19.8$ . This

result is confirmed by ordinary experience with the stock market. On any given day a stock value may easily move 3 to 5%, whereas the expected change is only about .05%. The daily mean is low compared to the daily variance.

## Mean Blur

We now show how this amplification effect makes the estimation of expected (or mean) rates nearly impossible.

Let us select a basic period length  $p$  (such as  $p = 1/12$  for a monthly period). We shall try to estimate the mean rate of return for this period. We assume that the statistical properties of the returns in each of the periods are identical, with mean value  $\bar{r}$  and standard deviation  $\sigma$ . We also assume that the individual returns are mutually uncorrelated. We wish to estimate the common mean value by using historical data.

Suppose that we have  $n$  samples of these period returns. The best estimate of the mean rate of return is obtained by averaging the samples. Hence,

$$\hat{r} = \frac{1}{n} \sum_{i=1}^n r_i. \quad (9.3)$$

The value of  $\hat{r}$  that we obtain this way is itself random. If we were to use a different set of  $n$  data points, we would obtain a different value of  $\hat{r}$ , even if the probabilistic character of the stock did not change (that is, if the true mean remained constant). However, the expected value of the estimate (9.3) is the true value  $\bar{r}$  since

$$E(\hat{r}) = E\left(\frac{1}{n} \sum_{i=1}^n r_i\right) = \bar{r}.$$

We want to calculate the standard deviation of the estimate  $\hat{r}$ , for it shows how accurate the estimate is likely to be. We have immediately

$$\sigma_{\hat{r}}^2 = E[(\hat{r} - \bar{r})^2] = E\left[\frac{1}{n} \sum_{i=1}^n (r_i - \bar{r})\right]^2 = \frac{1}{n} \sigma^2.$$

Hence,

$$\sigma_{\hat{r}} = \frac{\sigma}{\sqrt{n}}. \quad (9.4)$$

This is the basic formula for the error in the estimate of the mean value.

Let us put a few numbers into the formula. We take the period length to be 1 month. For the numbers used earlier, the monthly values are  $\bar{r} = 1\%$  and  $\sigma = 4.33\%$ . If we use 12 months of data, we obtain  $\sigma_{\hat{r}} = 4.33\%/\sqrt{12} = 1.25\%$ . Hence the standard deviation of the estimated mean is larger than the mean itself. If, using 1 year of data, we find  $\hat{r} = 1\%$ , we are only able to say, roughly, "the mean is 1% plus or minus 1.25%." This is not a good estimate. If we use 4 years of data, we cut this standard

deviation down by a factor of only 2—which is still poor. In order to get a good estimate, we need a standard deviation of about one-tenth of the mean value itself. This would require  $n = (43.3)^2 = 1,875$ , or about 156 years of data. However, the mean values are not likely to be constant over that length of time, and hence the estimation procedure is not really improved by much.

This is the historical blur problem for the measurement of  $\bar{r}$ . It is basically *impossible* to measure  $\bar{r}$  to within workable accuracy using historical data. Furthermore, the problem cannot be improved much by changing the period length. If longer periods are used, each sample is more reliable, but fewer independent samples are obtained in any year. Conversely, if smaller periods are used, more samples are available, but each is worse in terms of the ratio of standard deviation to mean value. (See Exercise 1.) The problem of mean blur is a fundamental difficulty.<sup>1</sup>

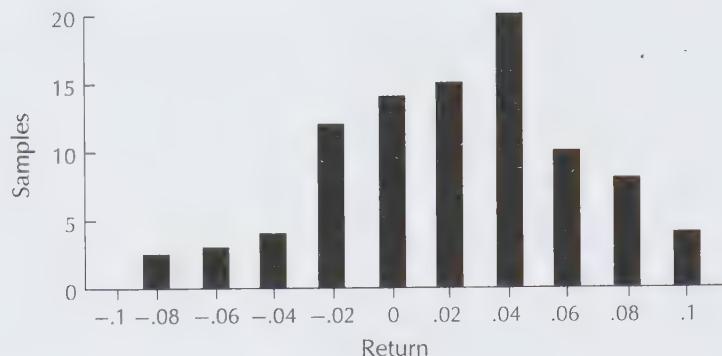
**Example 9.1 (A statistical try)** We simulated 8 years of monthly rates of return of a stock that had a monthly mean of 1% and a monthly standard deviation of 4.33%, corresponding approximately to yearly values of 12% and 15%, respectively. Random monthly returns were generated using a normal distribution with these parameters, and these returns are shown in the upper portion of Table 9.1. The sample means were calculated each year for the entire 8-year period. The sample means for each year are indicated below the monthly returns for that year. The sample standard deviation

**TABLE 9.1**  
**MONTHLY RATES OF RETURN AND ESTIMATION OF MEAN (EXPRESSED AS PERCENT)**

	Year of return								Overall
	1	2	3	4	5	6	7	8	
Jan	-8.65	2.61	6.39	-4.52	1.28	4.49	-1.44	3.30	
Feb	8.61	-2.38	-1.22	2.30	.14	7.58	-4.34	3.75	
Mar	5.50	-3.28	1.12	-3.96	-2.63	5.02	1.24	3.95	
Apr	2.04	7.45	3.69	-.84	3.15	-.51	8.92	-3.13	
May	7.51	7.96	.28	.35	-.47	-.19	-.46	-.31	
Jun	-2.50	-9.37	3.61	6.96	7.04	1.18	8.28	-.89	
Jul	2.28	-7.27	-1.45	4.23	3.68	1.61	-5.33	-6.39	
Aug	1.85	-5.30	6.83	.21	2.74	2.62	-1.01	-.60	
Sep	5.86	5.69	2.32	.14	-2.08	-2.32	3.77	-.76	
Oct	1.37	5.24	-3.79	-6.48	1.73	-3.08	4.18	1.92	
Nov	3.17	2.94	-.52	-1.11	6.18	5.42	-2.27	-3.97	
Dec	9.23	1.94	2.77	2.86	.38	2.93	4.91	5.18	
Mean	3.02	.52	1.67	.01	1.76	2.06	1.37	.17	1.32
$\sigma$	5.01	5.88	3.21	3.81	2.98	3.24	4.66	3.55	4.12

Each column represents a year of randomly generated returns. The true mean values are all 1%, but the estimates deviate significantly from this value.

<sup>1</sup> Some improvement can be made by estimating the means for several stocks simultaneously, as discussed in Section 9.4



**FIGURE 9.2 Histogram of monthly returns.** The distribution is too broad to pin down the true mean of .01 to within a small fraction of its value.

is also indicated. (Note that the sample standard deviations are also estimates—the accuracy of these is discussed in the following subsection.) Note also how the individual yearly estimates of the mean, as determined by the sample averages, jump around quite a bit from year to year. From this analysis we expect these estimates to have a standard deviation of 1.25%, and the results appear to be consistent with that. Even the 8-year estimate is quite far from the true value. We certainly should hesitate to use these estimates in a mean–variance optimization problem.

A histogram of the individual monthly returns is shown in Figure 9.2. Note that the standard deviation of the samples is large compared to the mean. One can see, visually, that it is impossible to determine an accurate estimate of the true mean from these samples. The mean value is too close to zero compared to the breadth of the distribution; hence one cannot pin down the estimate to within a small fraction of its actual value.

## 9.2 Estimation of Other Parameters

Estimates of other parameters from historical data are also subject to error. In some cases the error level is tolerable and in others it is not. In any event it is important to recognize the presence of errors and to determine their rough magnitudes—otherwise one might propose elaborate but fundamentally flawed procedures for portfolio construction.

### Estimation of $\sigma$

The blurring effect is not nearly as strong for the estimation of variances and covariances as it is for the mean. Suppose again that we have  $n$  samples of period rates of

return  $r_1, r_2, \dots, r_n$ . We calculate the sample mean

$$\hat{r} = \frac{1}{n} \sum_{i=1}^n r_i$$

and the sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (r_i - \hat{r})^2.$$

The use of  $n-1$  in the denominator instead of  $n$  compensates for the fact that  $\hat{r}$  is used instead of the true (but unknown)  $\bar{r}$ . It then follows that  $E(s^2) = \sigma^2$ . (See Exercise 4, Chapter 8.) Hence  $s^2$  provides an unbiased estimate of the variance.

The accuracy of the estimate  $s^2$  is given by its variance (or its standard deviation). It can be shown that if the original samples are normally distributed, the variance of  $s^2$  is

$$\text{var}(s^2) = \frac{2\sigma^4}{n-1}$$

or, equivalently,

$$\text{stddev}(s^2) = \frac{\sqrt{2}\sigma^2}{\sqrt{n-1}}.$$

This shows that the standard deviation of the variance is the fraction  $\sqrt{2/(n-1)}$  times the true variance, and hence the relative error in the estimate of  $\sigma^2$  is not too extreme if  $n$  is reasonably large.

**Example 9.2 (One year of data)** Suppose we again use a period length of 1 month. Using 12 months of data, we obtain  $\text{stddev}(s^2) = \sigma^2/2.35$ , which is already less than half of the value of  $\sigma^2$  itself. Hence the variance can be estimated with reasonable accuracy using about 1 year of historical data.

This conclusion is validated by the experiment shown in Table 9.1. The yearly estimates of  $\sigma$  shown in the bottom row are all reasonably close to the true value of 4.33% (certainly they are much better than the estimates of  $\bar{r}$ ), and the full 8-year estimate is really quite good.

## a Blur

The blur phenomenon applies to the parameters of a factor model, but mainly to the determination of  $a$ . In fact the presence of a blur can be deduced from the mean-blur phenomenon, but we omit the (somewhat complicated) details.

The same is true of  $\alpha$  for the security market line. It cannot be reliably estimated. On the other hand, the relative error in estimating  $\beta$  is somewhat better.

## 9.3 The Effect of Estimation Errors

Portfolio construction often relies on estimates of asset parameters, which, as we know, are in practice somewhat inaccurate, especially estimates of average rates of return but also, to a lesser degree, estimates of variances and covariances. It is prudent, then, to examine the likely impact that the use of incorrect parameter values may have on the quality of portfolios and to see if it is possible in fact to mitigate the negative consequences of those errors.

Ideally, we would like to quantify analytically the relation between parameter errors and their influence on portfolio performance; unfortunately, comprehensive analysis is available only for simple cases. However, some general statements and principles apply, and of course useful information can be obtained by simulation.

According to several such simulation studies, it is known that the negative impact of estimation errors on portfolio quality is significantly larger for expected-return errors than for variance and covariance errors. We also know from our previous discussion that it is very difficult to obtain reliable estimates of expected returns. Thus, expected-return errors are *both* greater and more damaging than estimation errors of variances and covariances. Therefore, it is reasonable to concentrate most analysis on the effect of expected-return errors. Accordingly, for simplicity, the remainder of this chapter assumes that estimates of return variances and covariances are reliable and that expected-return estimates are to some degree unreliable.

There are a number of ways that portfolios are negatively impacted by estimation errors:

**Condition number** The condition number of a matrix measures the degree of its singularity. If a matrix is close to being singular, it will have a high condition number, and the result of its operation (or that of its inverse) on a vector can be very sensitive to the values of vector components. In the case of portfolio design it is the covariance matrix  $\mathbf{V} = [\sigma_{ij}]$  that is of concern. If there is no risk-free asset, portfolios on the efficient frontier have weights  $w_j$ , determined from equation (6.5a) as

$$\sum_{j=1}^n \sigma_{ij} w_j - \lambda \bar{r}_i - \mu = 0$$

for various values of  $\lambda$  and  $\mu$ . It is clear that estimation errors of the  $\bar{r}_i$ 's might be magnified by a poorly conditioned covariance matrix  $\mathbf{V}$  and thus produce poor values for portfolio weights. A high condition number can frequently be traced to the fact that  $\mathbf{V}$  may be of large dimension, representing a large number of perhaps-similar assets.

**Leverage** A portfolio that has some weights greater than 1 and others less than zero—effectively buying the long positions with money obtained from the short positions—are particularly sensitive to estimation errors. For example, in the simplest case of just two assets, assume that the weights are  $w_1 = -5$  and  $w_2 = 6$ . The overall expected return is  $-5\bar{r}_1 + 6\bar{r}_2$ , which is highly sensitive to the difference  $\bar{r}_2 - \bar{r}_1$ . If this difference is small, errors can

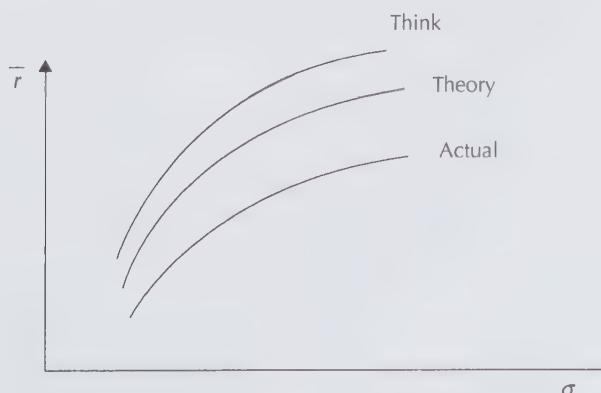
be greatly magnified. Even if there are no assets with weights greater than 1, some assets may have negative weights, and these serve to generate cash that can be devoted to other assets. This is similar to actual leverage and can have the same consequence of making the portfolio highly sensitive to the expected returns of a few assets.

**Large weights** If a portfolio includes assets with large weights, the overall return is highly dependent on the average rate of those assets. It is beneficial to diversify, for that not only decreases the market risk but also decreases the sensitivity of the expected overall return on the average return of a few assets.

## Three Views

Estimation errors cause us to produce portfolio designs that are not optimal. To study the effect of these errors, it is useful to consider three viewpoints—three different measures of design performance. Suppose we wish to construct the efficient frontier—obtained by optimizing expected return for various values of standard deviation. The results can be interpreted in three different ways. There is a curve of what we “think” is the efficient frontier, based on optimization using the estimates that we have used. Second is the curve that “theory” would predict if the true parameter values were used. Third, there is the “actual” curve that our portfolios would trace out using weights that we compute based on the estimates of the expected returns. Researchers have carried out extensive simulations that record these viewpoints. The resulting curves *on average* look more or less like those of Figure 9.3.

The “think” curve is what we calculate when optimizing the portfolio based on inaccurate estimates. Often this curve is an optimistic one, producing results that seem superior to those of “theory.” This can happen if the expected-return estimate of a particular asset is greater than the true value, for it will seem desirable to invest more



**FIGURE 9.3 Three ways to study the effect of estimation errors.** (1) What we “think” we will get, (2) what “theory” would predict assuming no estimation error, and (3) the “actual” result.

heavily in this asset, leading to an excessively positive portfolio result. However, it can go the other way as well, for a low estimate may lead to the downgrading of a good asset, causing us to invest less in this asset than would be optimal. However, as shown later, the “think” result will, on average, seem more favorable than that of “theory,” which uses true values.

Sadly, even though the performance we think we will get may be better than that determined by theory, in fact the portfolios we construct will always be worse than what theory predicts.

Let us systematically study the relations between the three different perspectives. Throughout this development, which spans the next few pages, it is convenient to define a separate multidimensional variable  $\mathbf{u}$  to denote the true expected-return vector  $\bar{\mathbf{r}}$ . This notation eliminates the need to use overlined variables and is commonly used in this type of analysis. With this convention, consider the problem

$$\begin{aligned} \max_{\mathbf{w}} \quad & \mathbf{w}^T \mathbf{u} \\ \text{subject to} \quad & \mathbf{w}^T \mathbf{V} \mathbf{w} \leq \sigma^2 \\ & \mathbf{w}^T \mathbf{1} = 1. \end{aligned} \tag{9.5}$$

As  $\sigma > 0$  is varied, the result traces out the efficient frontier of expected portfolio return versus its standard deviation  $\sigma$ . This is the “theory” curve.

Now consider the “think” case. We assume that the estimation process produces an unbiased estimate of  $\mathbf{u}$  of the form  $\mathbf{u} + \mathbf{e}$ , where  $E[\mathbf{e}] = \mathbf{0}$ . This estimate is used to maximize the objective  $\mathbf{w}^T(\mathbf{u} + \mathbf{e})$ . In particular, the “think” method solves

$$\begin{aligned} \max_{\mathbf{w}} \quad & \mathbf{w}^T(\mathbf{u} + \mathbf{e}) \\ \text{subject to} \quad & \mathbf{w}^T \mathbf{V} \mathbf{w} \leq \sigma^2 \\ & \mathbf{w}^T \mathbf{1} = 1 \end{aligned} \tag{9.6}$$

using the value of  $\mathbf{u} + \mathbf{e}$  that is observed, and we denote  $\mathbf{w}_{\mathbf{u}+\mathbf{e}}$  as the optimal weight vector corresponding to this perceived objective. The result is a random value that depends on  $\mathbf{e}$ .

The expected performance we anticipate at this “think” stage is

$$E\left[\max_{\mathbf{w}} \mathbf{w}^T(\mathbf{u} + \mathbf{e})\right], \tag{9.7}$$

where  $\mathbf{w}$  is subject to the two constraints for equation (9.6). On the other hand, the “theory” version can be expressed by reversing the order of maximization and expectation in equation (9.7), giving

$$\max_{\mathbf{w}} \{E[\mathbf{w}^T(\mathbf{u} + \mathbf{e})]\} = \max_{\mathbf{w}} \mathbf{w}^T \mathbf{u} \tag{9.8}$$

(with the same constraints). The  $\mathbf{e}$  disappears, since its expected value is  $\mathbf{0}$ . It can be shown<sup>2</sup> that the value of equation (9.7) is always greater than or equal to the value

---

<sup>2</sup> The objective function is linear in the random variable  $\mathbf{s} = \mathbf{u} + \mathbf{e}$ . It follows that  $\max_{\mathbf{w}} \mathbf{w}^T \mathbf{s}$  (subject to  $\mathbf{w}$  constraints) is a convex function of  $\mathbf{s}$ , say,  $f(\mathbf{s})$ . Jensen’s inequality then states that  $E[f(\mathbf{s})] \geq f(E(\mathbf{s}))$ .

of the left side of equation (9.8). Hence, on average, the “think” curve is higher than the “theory” curve.

Next consider the “actual” curve. The actual expected return is  $\mathbf{w}^T \mathbf{u}$ , but it uses the  $\mathbf{w}_{\mathbf{u}+\mathbf{e}}$  chosen in the design phase, which is optimal for maximizing  $\mathbf{w}^T(\mathbf{u} + \mathbf{e})$  but not for maximizing  $\mathbf{w}^T \mathbf{u}$ . The “theory” case produces the  $\mathbf{w}$  that is optimal for  $\mathbf{w}^T \mathbf{u}$ . Thus, since weights chosen at the “think” stage are not optimal for “theory,” the “actual” (achieved) value is always worse than the one that would be obtained in “theory.”

Summarizing the relations between the three views, we can write symbolically

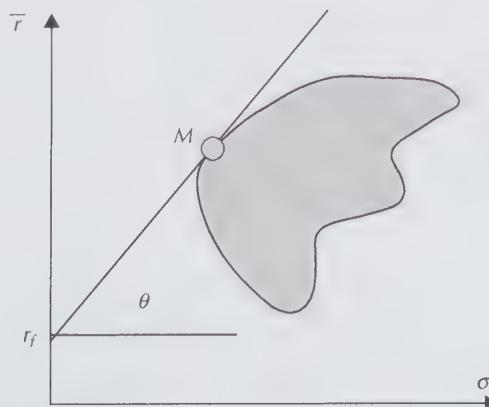
$$\text{Average Think} = E \text{ Max} > \text{Max } E = \text{Theory} > \text{Actual}.$$

Another interesting conclusion can be drawn. Notice that “theory” is always greater than “actual,” even if the sign of  $\mathbf{e}$  is reversed. It follows that near  $\mathbf{e} = \mathbf{0}$  the difference between “theory” and “actual” does not contain a linear component. The difference must be of second order in  $\mathbf{e}$ , for if it were linear, both positive and negative differences would occur near  $\mathbf{e} = \mathbf{0}$ .

The same argument applies to the relation between “average think” and “theory.” Since their difference is always positive regardless of the sign of  $\mathbf{e}$ , the difference is of second order in the estimation error  $\mathbf{e}$ .

## Maximum Tangent

In the context of mean–variance finance, the most important portfolio is the Markowitz (or market) portfolio. In the expected return versus standard deviation diagram, the line of maximum slope from the risk-free asset to the efficient set contains the Markowitz portfolio. See Figure 9.4. This line, at an angle  $\theta$ , is the capital market line; the slope is the price of risk. Because of its central role, it makes sense to investigate the sensitivity of that angle to estimation errors.



**FIGURE 9.4** The angle  $\theta$  of maximum tangent is defined by the Markowitz portfolio.

For simplicity, assume that the risk-free rate is zero. The basic problem for constructing a portfolio with maximum angle tangent is

$$\begin{aligned} \text{maximize}_w \quad & \frac{\mathbf{w}^T(\mathbf{u} + \mathbf{e})}{\sqrt{\mathbf{w}^T \mathbf{V} \mathbf{w}}} \\ \text{subject to} \quad & \mathbf{w}^T \mathbf{1} = 1, \end{aligned} \quad (9.9)$$

where  $\mathbf{V}$  is the covariance matrix of the returns. We assume that  $\mathbf{V}$  is known exactly. The actual expected returns are defined by the vector  $\mathbf{u}$ , but only the estimate  $\mathbf{u} + \mathbf{e}$  is available, where  $\mathbf{e}$  has expected value zero. This is the “think” problem for maximizing the tangent of the angle  $\theta$  under the assumption that the observed estimates are correct.

As before, this objective is linear in the random variable  $\mathbf{s} = \mathbf{u} + \mathbf{e}$ . Hence, the relations between the various viewpoints are the same

**Average Think > Theory > Actual.**

**Example 9.3 (Various angles)** Consider the case of two correlated risky assets and a risk-free asset. Without loss of generality, as before, we take  $r_f = 0$ . The other parameters are  $u_1 = 0.11$ ,  $u_2 = 0.06$ ,  $\sigma_1 = \sigma_2 = 0.20$ ,  $\rho = 0.7$ . From the general solution for two risky assets, equation (6.11), we find the optimal weights to be

$$w_1 = \frac{u_1 - u_2 \rho}{(u_1 + u_2)(1 - \rho)} = \frac{0.11 - 0.06 \times .7}{(0.11 + 0.06).3} = \frac{.068}{.051} = 1.333 \quad (9.10)$$

$$w_2 = 1 - w_1 = -.333. \quad (9.11)$$

These weights lead to an overall portfolio return of  $\bar{r} = w_1 u_1 + w_2 u_2 = 0.1267$ . The corresponding total variance is  $\sigma^2 = 0.04(w_1^2 + 2\rho w_1 w_2 + w_2^2) = 0.0507$ . Finally, the tangent of the critical angle is  $\tan(\theta) = \bar{r}/\sigma = 0.5627$ , which corresponds to an angle  $\theta$  of  $29.37^\circ$ . See the column labeled “Theory” in Table 9.2.

Next, we introduce an estimation error  $e = 0.02$ , which (in two separate cases) is either added to or subtracted from  $u_1$  (and hence has expected value 0). The corresponding results for “think” are in the “Think” columns of the table, and the results for “actual” are shown in the “Actual” columns. As predicted, the inequality relations between the three views is confirmed. Also, both the value of “actual” and the average value of “think” deviate from “theory” only by second order.

The relations in the example are somewhat encouraging. The two cases of “think” differ by first order, which means that the portfolio design is very much influenced by estimation error. Nevertheless, the resulting “actual” is always quite close to “theory.” In a sense, the system is somewhat immune to deviations that occur at the design stage. Furthermore, although the “think” error can be substantial, it is on average very close to “theory.” Perhaps this means that errors do not cause the havoc that we might have expected. Unfortunately, there is more to the story.

**TABLE 9.2  
THREE VIEWS**

	Think +	Think -	Theory	Actual +	Actual -
u1	0.11	0.11	0.11	0.11	0.11
u2	0.06	0.06	0.06	0.06	0.06
e	0.02	-0.02	0	0.02	-0.02
rho	0.7	0.7	0.7	0.7	0.7
w1	1.5439	1.0667	1.3333	1.5439	1.0667
w2	-0.5439	-0.0667	-0.3333	-0.5439	-0.0667
rbar	0.1681	0.0920	0.1267	0.1372	0.1133
sigma^2	0.0602	0.0417	0.0507	0.0602	0.0417
Tangent	0.6853	0.4505	0.5627	0.5594	0.5550
Delta	0.1225	-0.1122	0.0000	-0.0033	-0.0078
Average Delta	0.0052			-0.0056	
Theta	34.42	24.25	29.37	29.22	29.03

This table shows the solution to the maximum-angle tangent for two correlated risky assets and a risk-free asset with rate of return 0. The base case is the "Theory" column. Next, an estimation error is applied to  $u_1$  and this error has value  $e > 0$  in one case and  $-e < 0$  in another. They are assumed to be equally likely, so the expected value of their consequence is found by averaging the two corresponding results. Delta is the difference of a case from the "theory" case. As predicted, the expected tangent (the "Think" case) is either better or worse than the theoretical result, depending on the sign of  $e$ , but on average it is better. The "Actual" result is always lower than the theoretical value.

**TABLE 9.3  
EXAMPLE WITH AN ERROR OF  $\pm 0.10$** 

	Think +	Think -	Theory	Actual +	Actual -
u1	0.11	0.11	0.11	0.11	0.11
u2	0.06	0.06	0.06	0.06	0.06
e	0.1	-0.1	0	0.1	-0.1
rho	0.7	0.7	0.7	0.7	0.7
w1	2.0741	-1.5238	1.3333	2.0741	-1.5238
w2	-1.0741	2.5238	-0.3333	-1.0741	2.5238
rbar	0.3711	0.1362	0.1267	0.1637	-0.0162
sigma^2	0.0935	0.1323	0.0507	0.0935	0.1323
Tangent	1.2139	0.3744	0.5627	0.5355	-0.0445
Delta	0.6512	-0.1883	0.0000	-0.0273	-0.6072
Average Delta	0.2314			-0.3173	
Theta	50.52	20.53	29.37	28.17	-2.55

**Example 9.4 (A large-error case)** In practice, the estimation error associated with an asset of volatility 20% is likely to be much greater than 0.02. Four years of data would reduce the estimation error to about  $0.2/\sqrt{4} = 0.1$ . So let us examine the situation again with  $e = 0.10$ . The results are shown in Table 9.3. The damage

is much more severe than before, and the average deviations of “think” and “actual” from “theory” are *not* small (being .2314 and -.3173, respectively).

## Compounding Effect

There is another complication. Unlike the randomness associated with market volatility, estimation error is not reduced by time diversification, especially if the estimates are determined from historical studies. Those errors are not independent from period to period. In fact they are often identical or close to identical, since they are formed from fixed histories (slightly updated as time progresses). The same error basically repeats every period. The overall standard deviation after  $N$  periods is  $N\sigma$ ; hence, the ratio of sigma to expected value remains  $\sigma/\bar{r}$ . It does not decrease with time. In general, when there are both types of uncertainty, market volatility can be reduced by diversification; estimation error cannot.

## 9.4 Conservative Approaches

Poor portfolio performance can be upsetting, even more so when it is due to inaccurate parameter estimates. The recognition that poor performance may be due to estimation errors motivates one to design conservative portfolios, less sensitive to estimation errors. To search for such conservative methods, we recall the causes of estimation-error magnification described in the first part of the previous section.

Leverage can be eliminated by requiring portfolio weights to be nonnegative. Indeed, this is a very common portfolio restriction. In practice, it also tends to produce portfolios with relatively few nonzero weights, which has two other advantages. First, the resulting portfolio is easier to manage. Second, concentrating on only a few securities tends to reduce the condition number of the (smaller) covariance matrix and thus further reduce error sensitivity. However, it is desirable that no single asset have a large weight.

Another technique is explicitly to discourage large weights by the incorporation of a penalty term added to the portfolio objective function. Often a quadratic penalty is used, as in the following formulation (in the form of minimum variance rather than maximum expected return):

$$\begin{aligned} \text{minimize}_{\mathbf{w}} \quad & \mathbf{w}^T \mathbf{V} \mathbf{w} + c \mathbf{w}^T \mathbf{P} \mathbf{w} \\ \text{subject to} \quad & \mathbf{w}^T \bar{\mathbf{r}} = \bar{r} \\ & \mathbf{w}^T \mathbf{1} = 1. \end{aligned} \tag{9.12}$$

Here  $\mathbf{P}$  is a positive-definite matrix that ensures that  $\mathbf{w}^T \mathbf{P} \mathbf{w} > 0$  for all  $\mathbf{w} \neq \mathbf{0}$ . (Often  $\mathbf{P} = \mathbf{I}$  is suggested.) The constant  $c$  is a positive penalty factor that can be varied to produce different portfolio designs. The presence of the penalty term tends to force the optimization with respect to  $\mathbf{w}$  to decrease the magnitude of components of  $\mathbf{w}$  while still allowing some shorting. However, even if the restriction of nonnegative weights is imposed, the penalty term may be used as well.

A more direct technique is to set an upper bound on the weights. For example, it may be required that no asset be included in the portfolio with weight above 5%.

The conservative approach can be formalized in terms of robustness. In this approach, a probabilistic description of the estimation error is defined. Then an explicit trade-off between guarding against likely errors versus maximizing the objective is carried out. (See the references.)

## Better Estimates\*

The best way to get portfolios that perform as expected is to have accurate estimates of expected returns. Although the estimate based on the average of historical returns is the best possible unbiased estimate in the sense of mean-square error, there are other estimation methods that are not unbiased but have superior performance in a quadratic sense.

In particular, suppose there are  $n$  assets with covariance matrix  $\mathbf{V}$ . The true value of the expected rate of return is  $\mathbf{u}$ . We denote the average of the vector of observed rates of return as  $\hat{\mathbf{u}}$ . This is the standard estimator of  $\mathbf{u}$ , and it is unbiased. There are estimators  $\mathbf{u}^0$  that are superior to  $\hat{\mathbf{u}}$  in the sense of lower expected value of the loss function  $(\mathbf{u} - \mathbf{u}^0)^T \mathbf{V}^{-1} (\mathbf{u} - \mathbf{u}^0)$ . An important class is **shrinkage estimators** that use a weighted combination of two estimators. For instance, the first estimate might be  $\hat{\mathbf{u}}$ , and the second could be one that tends to bring the components close to each other. The estimate  $u_0 \mathbf{1}$  for some constant  $u_0$  is an example of a shrinkage estimate, since it makes all components equal. In fact, although some constants are better than others, *any*  $u_0$  will work as part of a two-part estimator.

A popular version is the **James–Stein Shrinkage estimator**,

$$\mathbf{u}_{JS} = (1 - w)\hat{\mathbf{u}} + w u_0 \mathbf{1}, \quad (9.13)$$

where

$$w = \min \left\{ 1, \frac{n - 2}{N(\hat{\mathbf{u}} - u_0 \mathbf{1})^T \mathbf{V}^{-1} (\hat{\mathbf{u}} - u_0 \mathbf{1})} \right\} \quad (9.14)$$

and the dimension of  $\mathbf{u}$  is  $n \geq 3$ . The fact that this estimator has lower expected quadratic loss than  $\hat{\mathbf{u}}$  for any  $u_0$  is known as **Stein's paradox**.

Another popular version of the shrinkage approach is applied in particular to the estimation of asset returns. It sets  $u_0$  equal to the estimated expected return of the minimum-variance portfolio (which is  $\mathbf{1}^T \mathbf{V}^{-1} \hat{\mathbf{u}} / \mathbf{1}^T \mathbf{V}^{-1} \mathbf{1}$ ) and

$$w = \frac{n + 2}{n + 2 + N(\hat{\mathbf{u}} - u_0 \mathbf{1})^T \mathbf{V}^{-1} (\hat{\mathbf{u}} - u_0 \mathbf{1})}. \quad (9.15)$$

It can easily be seen that  $0 < w < 1$ . In practice, shrinkage estimators applied to stock data tend to produce estimates of mean returns about 2–3% lower than the standard estimate.

Let us return to the discussion of the three viewpoints: think, theory, and actual. The most important measure is the difference between “think” and “actual.” This difference can be relatively large, and it is difficult to estimate its magnitude. On the other hand, although it is not nearly as useful, we know that the difference between “average think” and “actual” is of second order. We can use this to say something about the improvement associated with better estimates. Let  $\Delta$  be the difference between “average think” and “actual.” We know that  $\Delta$  is quadratic with respect to the magnitude of estimation errors; that is,  $\Delta = ce^2$  for some  $c > 0$  and some measure  $e$  of error magnitude. Suppose that our estimate is improved so that the error magnitude reduces to  $\alpha e$  for  $0 < \alpha < 1$ . We see that  $\Delta(\alpha) = \alpha^2 \Delta(1)$ , which is much better than what might have been expected. For instance, a 10% improvement in estimation (with  $\alpha = .90$ ) will ideally lead to a reduction in “average think” – “actual” of .81. In other words, a 10% improvement in estimation should lead to almost a 20% improvement in the difference.

This same conclusion applies to the difference between “theory” and “actual,” since it too is of second order.

**Example 9.5 (Earlier case)** In Table 9.3 with  $e = 0.1$ , if  $e$  is changed by 10% from 0.1 to 0.09, the resulting difference between “average think” and “theory” changes from 0.2314 to 0.1790, which is a reduction of 22.2%. Similarly, the resulting difference between average “actual” and “theory” changes from –0.3173 to –0.2690, a reduction of 15.2%.

## 9.5 Tilting Away From Equilibrium\*

Better estimates of mean returns can be obtained if there is information regarding the future prospects of the stock available that supplements the information contained in the historical record. Such information can be obtained in a variety of ways, including: (1) from detailed fundamental analyses of the firm, including an analysis of its future projects, its management, its financial condition, its competition, and the projected market for its products or services, (2) as a composite of other analysts’ conclusions, or (3) from intuition and hunches based on news reports and personal experience. Such information can be systematically combined with the estimates derived from historical data to develop superior estimates.

One potential source is obtained by using the CAPM in a reverse fashion. It determines the expected rates of return that would be required to produce the market portfolio. That is, a set of expected rates of return is found, which, when used as the rates in the mean–variance problem, will lead to the market portfolio as the solution. Let us see how that works.

The required CAPM rates are given by the CAPM formula; namely,

$$\bar{r}_i^e = r_f + \beta_i(\bar{r}_M - r_f).$$

We have added the superscript  $e$  to emphasize that this is the value of  $\bar{r}_i$  obtained through the equilibrium argument. Note that this value of  $\bar{r}_i^e$  is fairly easy to obtain. It is only necessary to estimate  $\beta_i$  (which can be estimated quite reliably) and  $\bar{r}_M$  (which is more difficult, but often a consensus view can be used). No equations need be solved.

The true expected rates of return are random variables that we cannot know with certainty. The equilibrium values computed before give us some information about these values, but these too are only estimates. We expect that these estimates each have some variance and they are correlated with each other. We therefore write the equation

$$\bar{r}_i = \bar{r}_i^e + \varepsilon_i$$

for each stock  $i$  to express the fact that the true value of  $\bar{r}_i$  is equal to the values obtained by the equilibrium argument plus some error. The error  $\varepsilon_i$  has zero mean.

Other information about expected rates of return can be expressed in a similar way. For example, to incorporate historical data on asset  $i$ , we might write an equation of the form  $\bar{r}_i = \bar{r}_i^h + e_i$ , where  $\bar{r}_i^h$  is the value of  $\bar{r}_i$  obtained from historical data and  $e_i$  has variance equal to that implied by the length of the historical record.

Likewise, we might include subjective information about the expected return, or information based on a careful analysis of the firm. In each case we also assign a variance to the estimate.

We can imagine building up the estimate in steps. We can start with the estimate based on the equilibrium expected returns. This will lead to the market portfolio as the solution to the Markowitz problem. As additional information is added, the solution will **tilt** away from that initial solution. The degree of departure, or tilt, will depend on the nature of the adjoined equations and the degree of confidence we have in them, as expressed by the variances and covariances of the error terms. This can be regarded as a version of shrinkage estimation since it combines two or more estimates according to some weights.

**Example 9.6 (A double use of data)** Refer to Example 8.2 and the data of Table 8.2. Most of the summary part of this table is repeated here in Table 9.4. The first row of the table gives the 10-year average returns. It is easy to calculate the corresponding CAPM estimates. For example, for stock 1 we have  $\bar{r}_1^e = 5.84 + .90(13.83 - 5.84) = 13.05$ . These estimates are clearly not equal to the historical averages.

To form new, combined, estimates, we assign a variance to each estimate. Since there are 10 years of data, it is appropriate to use (9.4) to write  $\sigma_i^h = \sigma_i / \sqrt{10}$  for the standard deviation of the error in the historical estimate of  $\bar{r}_i$ . For stock 1, this is  $\sigma_1^h = \sqrt{90.28/10} = 3.00$ .

To assign error magnitudes to the CAPM estimates, we notice that these estimates are based on our estimates of  $r_f$ ,  $\beta_i$ , and  $\bar{r}_M$ . Let us ignore all errors except that contained in  $\bar{r}_M$ . The standard deviation of the error in  $\bar{r}_1^e$  is thus  $\beta_1 \times \sigma_M / \sqrt{10} = .90 \sqrt{72.12/10} = 2.42$ .

**TABLE 9.4**  
**DATA FOR TILTING**

	Stock 1	Stock 2	Stock 3	Stock 4	Market	Riskless
aver	15.00	14.34	10.90	15.09	13.83	5.84
var	90.28	107.24	162.19	68.27	72.12	
cov	65.08	73.62	100.78	48.99	72.12	
$\beta$	.90	1.02	1.40	.68	1.00	
CAPM	13.05	14.00	17.01	11.27		
tilt	13.82	14.14	14.17	12.52		

The historical average returns are not equal to the average returns predicted by CAPM. Both estimates have errors, but they can be combined to form new estimates, called *tilt*.

For stock 1, if we treat these two estimates of  $\bar{r}_1$ , the historical and the CAPM (equilibrium) estimates, as independent, then they are best combined by<sup>3</sup>

$$\bar{r}_1 = \left[ \frac{\bar{r}_1^h}{(3.00)^2} + \frac{\bar{r}_1^e}{(2.42)^2} \right] \left[ \frac{1}{(3.00)^2} + \frac{1}{(2.42)^2} \right]^{-1} = 13.82.$$

(See Exercise 9.) The new estimates for the other stocks are found in a similar fashion.

## 9.6 Summary

It is tempting to assume that the parameter values necessary to implement mean-variance theory—namely, the expected returns, variances, and covariances for the Markowitz formulation, or the  $a_i$ 's and  $b_{ij}$ 's for a factor model—can be estimated from historical returns. Although some parameters can be estimated this way, others cannot. In particular, for stocks the variances and covariances can be estimated to within reasonable accuracy using about 1 year of weekly or daily returns data. However, the expected returns (the means) are subject to a blurring phenomenon and cannot be estimated to within workable accuracy, even if a record of 10 years of returns is available. This blurring phenomenon applies to the  $a$  coefficients in a factor model as well.

Still, investors should make the best of what they have. Every portfolio design should account for such errors and provide some means for minimizing their impact. Especially important are errors in the estimation of asset expected returns, since these are usually relatively large compared to errors in the estimation of variances and covariances, and they also seem to have a greater deleterious effect on portfolio results. It is useful to analyze the impact of estimation errors in terms of three viewpoints: (1) what a designer thinks, that is, what portfolio average return the designer predicts

<sup>3</sup> These two estimates are not really independent since the historical market return is based in part on the historical return of stock 1. Furthermore, the CAPM errors of different stocks are highly correlated since they all depend on the market. We ignore these correlations for the sake of simplicity.

based on available estimates; (2) what theory would predict if estimates were accurate; and (3) what expected return will actually be achieved. It can be shown that on average the designer's expectation will be greater than the expectation computed in theory, and that in turn will be greater than what is actually achieved. If estimation errors are small, the differences between these views are even smaller. But, unfortunately, errors are not always small. Even more troublesome is the fact that estimation error is not diversifiable, since the same (or a similar) error occurs in each period.

The effect of estimation errors can to some extent be mitigated by conservative design, such as requiring that assets weights be nonnegative and not large.

Finally, it is possible, by estimating the expected returns of several assets simultaneously, to achieve better estimates than the ones that would be obtained by estimating the expected returns of each asset separately, using the method of shrinkage estimation or by tilting away from equilibrium.

## Exercises

1. (Are more data helpful?  $\diamond$ ) Suppose a stock's rate of return has annual mean and variance of  $\bar{r}$  and  $\sigma^2$ . To estimate these quantities, we divide 1 year into  $n$  equal periods and record the return for each period. Let  $\bar{r}_n$  and  $\sigma_n^2$  be the mean and the variance for the rate of return for each period. Specifically, assume that  $\bar{r}_n = \bar{r}/n$  and  $\sigma_n^2 = \sigma^2/n$ . If  $\hat{\bar{r}}_n$  and  $\hat{\sigma}_n^2$  are the estimates of these, then  $\hat{\bar{r}} = n\hat{\bar{r}}_n$  and  $\hat{\sigma}^2 = n\hat{\sigma}_n^2$ . Let  $\sigma(\hat{\bar{r}})$  and  $\sigma(\hat{\sigma}^2)$  be the standard deviations of these estimates.

- (a) Show that  $\sigma(\hat{\bar{r}})$  is independent of  $n$ .
- (b) Show how  $\sigma(\hat{\sigma}^2)$  depends on  $n$ . (Assume the returns are normal random variables.)

Answer the question posed as the title to this exercise.

2. (A record) A record of annual percentage rates of return of the stock  $S$  is shown in Table 9.5.

**TABLE 9.5**  
**RECORD OF RATES OF RETURN**

Month	Percent rate of return	Month	Percent rate of return
1	1.0	13	4.2
2	.5	14	4.5
3	4.2	15	-2.5
4	-2.7	16	2.1
5	-2.0	17	-1.7
6	3.5	18	3.7
7	-3.1	19	3.2
8	4.1	20	-2.4
9	1.7	21	2.7
10	.1	22	2.9
11	-2.4	23	-1.9
12	3.2	24	1.1

- (a) Estimate the arithmetic mean rate of return, expressed in percent per year.  
 (b) Estimate the arithmetic standard deviation of these returns, again as percent per year.  
 (c) Estimate the accuracy of the estimates found in parts (a) and (b).  
 (d) How do you think the answers to (c) would change if you had 2 years of weekly data instead of monthly data? (See Exercise 1.)
3. (Clever, but no cigar  $\diamond$ ) Gavin Jones figured out a clever way to get 24 samples of monthly returns in just over one year instead of only 12 samples; he takes overlapping samples; that is, the first sample covers Jan. 1 to Feb. 1, and the second sample covers Jan. 15 to Feb. 15, and so forth. He figures that the error in his estimate of  $\bar{r}$ , the mean monthly return, will be reduced by this method. Analyze Gavin's idea. How does the variance of his estimate compare with that of the usual method of using 12 nonoverlapping monthly returns?
4. (Tangent errors) For Example 9.3, represent the uncertainty in the expected returns of the two stocks in the form of  $\mathbf{u} + \mathbf{e}$ , where  $E[\mathbf{e}] = \mathbf{0}$  with covariance matrix of  $\mathbf{e}$  equal to  $\mathbf{Q}$ .
5. (Volatility sensitivity) In Table 9.2, suppose the variance of the two assets are incorrectly estimated to be .044 instead of the actual .04 (which is a 10% error). Find the "think" value for the tangent of the angle of the capital market line.
6. (A plot) Plot the value of "actual"—"theory" for Example 9.3 for  $e$  in the range  $0 \leq e \leq .20$ . Explain the rather dramatic character of the plot.
7. (Nonnegativity constraints) Suppose there are three assets with expected rates of return  $r_1 = 0.10$ ,  $r_2 = 0.07$ ,  $r_3 = 0.14$ . Each has a standard deviation of .20, and all pairs have correlation coefficient  $\rho = .4$ . There is a constraint that the variance of the portfolio must be less than 0.038.
- (a) Find the "theory" value.  
 (b) Find the value when the weights are constrained to be nonnegative.
8. (Minimum-variance formulation\*) Consider this problem:
- $$\begin{aligned} &\text{minimize}_{\mathbf{w}} \quad \mathbf{w}^T \mathbf{V} \mathbf{w} \\ &\text{subject to} \quad \mathbf{w}^T (\mathbf{u} + \mathbf{e}) \geq \bar{r} \\ &\quad \mathbf{w}^T \mathbf{1} = 1. \end{aligned}$$
- Let "think" be the value with  $\mathbf{e}$ , and let "theory" be the value with  $\mathbf{e} = \mathbf{0}$ . Now assume that  $\mathbf{e}$  is random with expected value zero. Show that the expected value of "think" is less than the value of "theory." [Hint: First show (or use) the fact that the value of "think" is a concave function of  $\mathbf{e}$ . Then use Jensen's inequality.]
9. (General tilting  $\diamond$ ) A general model for information about expected returns can be expressed in vector-matrix form as

$$\mathbf{p} = \mathbf{P}\bar{\mathbf{r}} + \mathbf{e}.$$

In the model  $\mathbf{P}$  is an  $m \times n$  matrix,  $\bar{\mathbf{r}}$  is an  $n$ -dimensional vector, and  $\mathbf{p}$  and  $\mathbf{e}$  are  $m$ -dimensional vectors. The vector  $\mathbf{p}$  is a set of observation values and  $\mathbf{e}$  is a vector of errors having zero mean. The error vector has a covariance matrix  $\mathbf{Q}$ . The best (minimum-variance) estimate of  $\bar{\mathbf{r}}$  is

$$\hat{\mathbf{r}} = (\mathbf{P}^T \mathbf{Q}^{-1} \mathbf{P})^{-1} \mathbf{P}^T \mathbf{Q}^{-1} \mathbf{p}. \quad (9.16)$$

- (a) Suppose there is a single asset and just one measurement of the form  $p = \bar{r} + e$ . Show that according to equation (9.16), we have  $\hat{r} = p$ .
- (b) Suppose there are two uncorrelated measurements with values  $p_1$  and  $p_2$ , having variances  $\sigma_1^2$  and  $\sigma_2^2$ . Show that

$$\hat{r} = \left( \frac{p_1}{\sigma_1^2} + \frac{p_2}{\sigma_2^2} \right) \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} \right)^{-1}.$$

- (c) Consider Example 9.6. There are measurements of the form

$$\hat{r}_1 = p_1 + e_1$$

$$\hat{r}_2 = p_2 + e_2$$

$$\hat{r}_3 = p_3 + e_3$$

$$\hat{r}_4 = p_4 + e_4$$

$$\hat{r}_1 = r_f + \beta_1 f_M$$

$$\hat{r}_2 = r_f + \beta_2 f_M$$

$$\hat{r}_3 = r_f + \beta_3 f_M$$

$$\hat{r}_4 = r_f + \beta_4 f_M,$$

where the  $e_i$ 's are uncorrelated but where  $\text{cov}(e_i, f_M) = .25\sigma_i^2$ . Using the data of the example, and assuming the  $\beta_i$ 's are known exactly, find the best estimates of the  $\bar{r}_i$ 's.  
[Note: You should only need to invert  $2 \times 2$  matrices.]

## References

The analysis of errors in the estimation of return parameters from historical data has long been available, but it is not widely emphasized. See [1] for an early good treatment. The effect of estimation errors on portfolios was studied by simulation and reported in [2]. The idea of comparing portfolio results from different viewpoints goes back to [3] and [4]. A good overview of robustness for portfolio design is [5]. An excellent presentation of robustness and shrinkage estimation is [7]. Also see [6]. The original references on shrinkage estimators are [9] and [10] and for asset problems in particular [11]. The tilting procedure was proposed by Black and Litterman [8].

1. Ingersoll, J. E. (1987), *Theory of Financial Decision Making*, Rowman and Littlefield, Savage, MD.
2. Chopra, V., and W. Ziemba (Winter 1992), "The Effect of Errors in Means, Variances, and Covariances on Optimal Portfolio Choice," *Journal of Portfolio Management*, 6–11.
3. Broadie, M. (1993), "Computing Efficient Frontiers Using Estimated Parameters," *Annals of Operations Research*, **45**, nos. 1–4, 21–58.
4. Ceria, S., and R. Stubbs (2006), "Incorporating Estimation Errors in Portfolio Selection: Robust Portfolio Construction," *Journal of Asset Management* **7**, no. 2, 109–127.
5. Shaw, D. (February 2008), "Robust Optimization: What Works and What Does Not," *Northfield News* (slides), [www.northinfo.com/documents/285.pdf](http://www.northinfo.com/documents/285.pdf)
6. Kuhn, D., P. Parpas, B. Rustem, and R. Fonseca, (2009), "Dynamic Mean–Variance Portfolio Analysis under Model Risk," *Journal of Computational Finance*, **12**, no. (4), 91–115.
7. Fabozzi, F., P. Kolm, D. Pachamanova, and S. Forcardi (2007), *Robust Portfolio Optimization and Management*, John Wiley & Sons, Hoboken, NJ.

8. Black, F., and R. Littleman (September/October 1992), "Global Portfolio Optimization," *Financial Analysts Journal*, 28–43.
9. Stein, C. (1995) "Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution," *Proceedings of the Third Berkeley Symposium on Mathematical Probability and Statistics, I*. Berkeley: University of California Press, 197–206.
10. James, W., and C. Stein (1961), "Estimation with Quadratic Loss," *Proceedings of the Fourth Berkeley Symposium on Mathematical Probability and Statistics, I*. Berkeley: University of California Press, 361–379.
11. Jorion, P. (Sept. 1986), "Bayes–Stein Estimation for Portfolio Analysis," *The Journal of Financial and Quantitative Analysis*, 21, no. 3, 279–292.

# 10

## RISK MEASURES

Prudent investors are concerned with risk, and they generally design portfolios on the basis of a comfortable trade-off between the risk of loss and the possibility of profit. This can be done explicitly by the use of a quantitative risk measure.

One popular measure of this trade-off is the Sharpe ratio of a portfolio—equal to the average excess return (above the risk-free rate) divided by the standard deviation of the return of the portfolio; that is,  $(\bar{r}_p - r_f)/\sigma_p$ . (See Chapter 7.) This single ratio is a simple yet practical measure of the quality of a portfolio, and it is in common use. In some cases, however, it is desirable to skew the trade-off heavily, seeking to focus on risk of loss alone, with the possibility of gain taking a back seat.

Banks and other financial intuitions, for instance, represent situations where attention paid to the risk of loss overwhelmingly dominates the anticipation of possible gain. The promised gain of direct deposits and certificates of deposit is essentially known exactly, but the chance for loss is not easily estimated. There is some chance, generally quite small, that the bank, which itself invests in risky assets, may lose a great deal of money and hence, through default, jeopardize the promised return to individuals.

Since it is not practical for every bank customer to assess a bank's risk carefully, the task has to a large extent been taken up by the government in the form of regulations, oversight, and insurance—in the United States by the Federal Deposit Insurance Corporation (FDIC). Bank regulators use standardized risk measures to assess soundness, expressed mainly in terms of the probability of a large loss. Financial soundness typically comes down to whether the bank has sufficient capital reserves to protect against the possibility of a major loss; indeed, regulators typically focus only on

potential loss, not on gain. That is the viewpoint of this chapter—a focus exclusively directed at quantifying the risk of loss.

Risk can be divided into a few categories. One of the most important, **market risk**, is the risk associated with assets that can be freely traded in the market, such as public stocks and bonds. A portfolio consisting of marketed assets can be evaluated every day, and hence a projected value for the next day can be estimated probabilistically. **Nonmarket risk** applies to assets that are held but are not directly priced by the market on a short-term basis. A nontraded loan is an example. Investment in a private firm is another. The lack of liquidity for such nonmarket investments makes them difficult to value. **Credit risk**, the risk of counterparty default on a loan or other transaction, is another major category of risk. It can be viewed as a combination of market risk and nonmarket risk, since default may be caused at least in part by general market conditions. The topic of credit risk is important; significant and expanding theory is devoted to it. It is treated in more detail in Chapter 17, but some of the underlying ideas spring from the material in this current chapter.

## 10.1 Value at Risk

Imagine a financial position that defines a random value  $X$  denoting the change in the value of the position at some given future time  $T$ . In general, the variable  $X$  may take on either positive or negative values, depending on its realization (outcome). For convenience we refer to the random variable  $X$  itself as a **position**. From a risk perspective, we may focus on the associated loss, which is  $-X$ .

Indeed, **value at risk** (abbreviated VaR) is motivated by the concern about loss. To define this notion, one specifies a **loss tolerance**  $h$  between zero and 1 and a companion **confidence level** equal to  $1 - h$ . For example, one might choose a loss tolerance of  $h = .05$  and a corresponding confidence level of  $1 - h = .95$  or, equivalently, a loss tolerance of 5% and a confidence level of 95%. For a particular position  $X$  and a given loss tolerance  $h$ , VaR is then the smallest number  $V$  such that the probability of a loss greater than  $V$  is no more than  $h$ . Mathematically<sup>1</sup>

$$\text{VaR}_h(X) = \min_h \left\{ V : P[-X > V] \leq h \right\}. \quad (10.1)$$

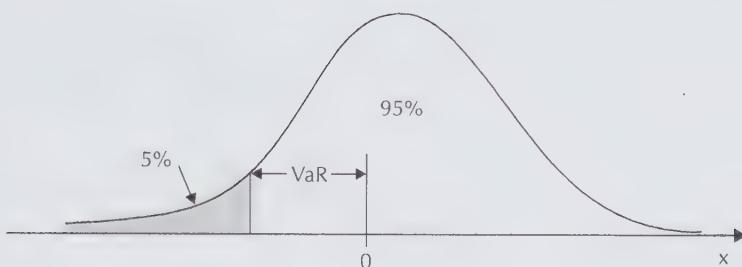
Equivalently, VaR is the smallest number  $V$  such that the probability of the loss being no more than  $V$  is greater than  $1 - h$ . That is,

$$\text{VaR}_h(X) = \min_h \left\{ V : P[-X \leq V] > 1 - h \right\}. \quad (10.2)$$

Value at risk is illustrated in Figure 10.1, which shows a hypothetical probability density of  $X$ . The total area under this curve is, of course, 1. The VaR is determined by the point at the left (lower) end, where the total area under the density above that point is equal to the confidence level, say, 95% of the total. Equivalently, it is the point

---

<sup>1</sup>  $P[\cdot]$  denotes the probability of whatever event is between the brackets. If the minimum does not exist, the least upper bound (or infimum), denoted  $\inf$ , is used.



**FIGURE 10.1 VaR as the lower range of probability for a 95% confidence level.** The probability density shown has an expected value that is somewhat positive, and hence the zero point is to the left of center. VaR is measured from that zero point to the critical quantile, defined by the loss tolerance.

where the total area below the VaR point is equal to the loss tolerance 5% of the total. The point on the axis that defines the upper boundary of the region under the density with the given area is termed a **quantile** or, specifically for an area of magnitude  $h$ , the  $h$ -quantile. Value at risk is therefore equal to the quantile at a specific loss tolerance, but changed in sign to measure loss rather than (negative) gain.

For the moment, we assume that the density function is a smooth, continuous function. In particular, we assume that there are no discrete probabilities, or “atoms” (points having positive probability).

Implicit in the definition of VaR is the time  $T$  at which  $X$  is realized. For monitoring of bank reserves, this is usually one or a few days into the future. For situations where there are no highly liquid markets, the period may be as long as a year. The confidence level is specified by the regulating agency, and it is often 95% or 99%, although other values may be used for the bank’s internal purposes.

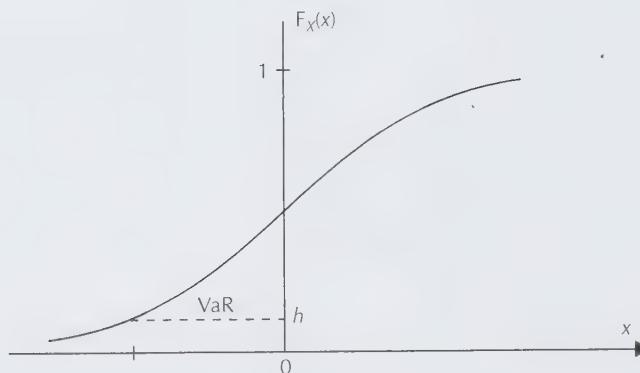
We use the notation  $\text{VaR}_h(X)$  for the value at risk of  $X$  at loss tolerance of  $h$  (equivalent to confidence level of  $1 - h$ ).

It is useful to consider the **cumulative probability distribution** of  $X$ , defined as  $F_X(x) = P[X \leq x]$ , where  $P$  denotes probability. See Figure 10.2. At any point  $x$ ,  $F_X(x)$  is the total area under the density to the left of  $x$ . The distribution is nondecreasing, moving from 0 to 1. However, there can be flat spots (if there is a range where  $X$  has zero probability) or jumps (as in the case of discrete probabilities). If a flat spot occurs with  $F_X(x) = h$  for an interval of  $x$ ’s, then  $\text{VaR}_h(X)$  is defined by essentially taking (the negative of) the rightmost point of the flat interval. This means that we take  $\text{VaR}_h(X)$  to be the *least* possible value among the contenders.<sup>2</sup>

If  $F_X(x)$  is continuous and strictly increasing, it is possible to define the inverse function  $x = F_X^{-1}(h)$ . In terms of this inverse we have

$$\text{VaR}_h(X) = -F_X^{-1}(h). \quad (10.3)$$

<sup>2</sup> Specifically, in all cases  $\text{VaR}_h(X) = -\inf\{m : F_X(m) > h\}$ .



**FIGURE 10.2 Cumulative probability distribution and its relation to VaR.**

A method to determine the inverse function  $F_X^{-1}(x)$  is illustrated in Figure 10.2. To find the value at a particular value  $x = h$ , one finds  $h$  on the vertical axis and moves horizontally to the point of intersection with the curve; the corresponding value on the horizontal axis is  $F_X^{-1}(h)$ . Note that for most distributions of concern, a small  $h$  produces a negative value and, correspondingly, a positive VaR.

## Properties of VaR

Value at risk satisfies some basic properties that help characterize its structure and clarify the degree to which it is a useful measure of risk. First, VaR is easily understood, and it intuitively reflects potential loss. Second, VaR depends only on the left tail of the density, which means that it focuses exclusively on loss without (perhaps optimistic) projections of gain. Third, value at risk is objective, in the sense that all assumptions are reduced to probabilities rather than being defined by broad categories, such as AA, A, BB etc. as typically provided by rating companies. Of course, in application one must determine carefully the probability density in the left tail, and this may be difficult.

## Capital Requirement

One of the most important properties of VaR is that it can be interpreted as the amount of risk-free payoff that must be added to a position  $X$  in order that the new VaR becomes zero. The VaR point is located at the point  $-\text{VaR}_h(X)$  on the axis of the density function or distribution function. See Figures 10.1 and 10.2. If an amount  $C > 0$  is added to  $X$ , the density graph will move to the right by an amount  $C$ . In particular, if one adds  $C = \text{VaR}_h(X)$ , then the new value at risk will be zero. Hence,  $\text{VaR}_h(X) = C$ , where  $C$  is the least capital payoff that must be added to  $X$  to ensure

that a loss greater than zero has probability no greater than  $h$ . In other words,<sup>3</sup> VaR is the answer to the question “How much capital payoff do I need to add so that a loss will occur with only probability  $h$ ? ”

## 10.2 Computation of Value at Risk

Determination of the numerical value of VaR requires knowledge or an estimate of the probability distribution of the position  $X$ , especially in the lower tail of the density, where large potential losses are represented. We restrict our attention here to the effects of market risk only, for which systematic appraisal is relatively straightforward, even if not easy.

### Model-Based Method

Often it is assumed that the probability distribution of  $X$  can be represented analytically. Likewise, the appropriate quantile can then be found analytically or by computational approximation.

A common assumption is that the risks of individual assets in a portfolio follow normal (that is, Gaussian) probability densities. Such densities are commonly used in statistical analysis, as discussed in detail in Appendix A and in Chapters 13 and 15. For our present purpose, only a few simple facts are needed. The normal probability density function<sup>4</sup> is represented by the familiar **bell-shaped curve** shown in Figure 10.3. It is completely specified by its mean value  $\mu$  and variance  $\sigma^2$  (or standard deviation  $\sigma$ ); these of course are the same two parameters frequently used in earlier chapters to characterize uncertain returns.

Given  $X$ , which is normal and has mean value  $\mu$  and variance  $\sigma^2$ , the associated VaR can be found by use of Proposition 10.1. (Readers unfamiliar with the distribution of the normal variable may safely skip the proof and consider the examples that follow.)

**Proposition 10.1 (VaR for normal distribution)** Suppose  $X$  follows a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ . Then

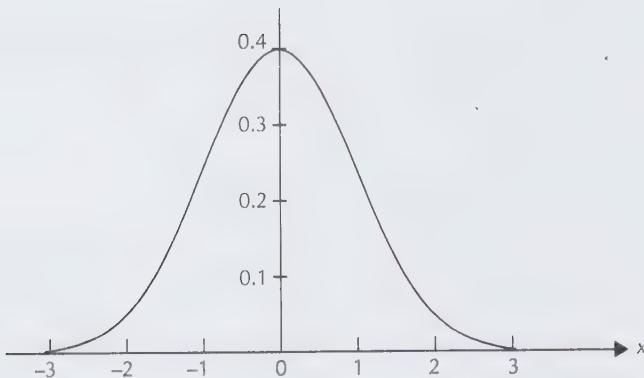
$$\text{VaR}_h(X) = -\sigma F_N^{-1}(h) - \mu, \quad (10.4)$$

where  $F_N$  is the cumulative probability distribution function of the standardized normal variable (with mean zero and standard deviation 1).

<sup>3</sup> We have (see Exercise 4.)  $\text{VaR} = \inf\{C : P[X + C < 0] \leq h\}$ .

<sup>4</sup> The normal density function of mean  $\mu$  and standard deviation  $\sigma$  is

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}.$$



**FIGURE 10.3 The standard normal probability density.** It has mean of zero and standard deviation of 1. (The tails go out to infinity in both directions.)

**Proof:** It is convenient to define a new normal random variable  $X^* = (X - \mu)/\sigma$ , which is standardized, in the sense that it has mean 0 and standard deviation 1. This standardized normal distribution is referred to as  $N(0, 1)$ , and the associated values of the distribution function<sup>5</sup> are available by use of standard tables or software packages. See Exercise 1, Chapter 15.

To find the value at risk for  $X$  we let  $m = -\text{VaR}_h(X)$ . Then, starting with  $P[X \leq m] = h$ , we change both sides of the inequality by subtracting  $\mu$  and dividing by  $\sigma$  to obtain

$$\begin{aligned} P[(X - \mu)/\sigma \leq (m - \mu)/\sigma] &= h \\ P[X^* \leq (m - \mu)/\sigma] &= h \\ F_N[(m - \mu)/\sigma] &= h. \end{aligned} \tag{10.5}$$

We then solve (10.5) for  $m$  and take the negative of the result. Hence

$$\begin{aligned} (m - \mu)/\sigma &= F_N^{-1}(h) \\ \text{VaR}_h(X) &= -\sigma F_N^{-1}(h) - \mu. \blacksquare \end{aligned}$$

Notice that for small values of  $h$  (indeed for any  $h$  less than .5), the value of  $F_N^{-1}(h)$  is negative. Hence, VaR increases as standard deviation increases.

**Example 10.1 (Small  $\mu$ )** Consider a highly liquid portfolio whose payoff is the normal random variable  $X$ . If  $T$  is small, such as a few days, the mean value of  $X$  is likely to be much smaller than the standard deviation. If we assume, as a limiting case, that  $\mu = 0$ , then the value at risk depends only on the variance. For example,

<sup>5</sup> The notation  $P[X^* \leq x]$  again denotes the probability of the event  $X^* \leq x$ . Thus  $F_N(x) = P[X^* \leq x]$ .

take  $h = .05$ . We can find the value at risk from equation (10.4). Using a table of the normal distribution we find  $-F_N^{-1}(.05) = 1.65$ . This means that  $\text{VaR}_{.05}(X) = 1.65\sigma$ .

**Example 10.2 (N-day VaR)** Suppose a pension fund invests its endowment of \$10 million in a position  $X_A$  that follows a normal distribution. The fund wishes to calculate the 10-day VaR of this position at a confidence level of  $1 - h = .95$ . For such calculations, it is common to express means and variances of the position on a per-day rather than a per-year basis as we have done for most analyses in other chapters. If the yearly values of these parameters are  $\mu_y$  and  $\sigma_y$ , then the corresponding daily values are  $\mu_d = \mu_y/250$  and  $\sigma_d = \sigma_y/\sqrt{250}$ . (See Chapter 9.) It is also common to assume in these calculations, as mentioned earlier, that the daily value of  $\mu$  is zero, since indeed in most cases it will be very small compared to the daily  $\sigma$ .

Assume that  $X_A$  has a volatility of 2% per day (which is about 30% per year), and suppose that the returns on different days are independent. Then for  $N = 10$  days, the volatility increases by the square root of 10. Hence we take  $\mu = 0$  and  $\sigma = .02 \times \sqrt{10} = 6.32\%$  for 10 days.

We readily find from equation (10.4) that

$$\begin{aligned}\text{VaR}_{.05}(X_A) &= \$10 \text{ million} \times [-.0632 \times F_N^{-1}(.05)] \\ &= \$10 \text{ million} \times .0632 \times 1.65 \\ &= \$1.0428 \text{ million.}\end{aligned}$$

**Example 10.3 (Diversification)** Suppose that the pension fund in the previous example also invests \$10 million in another position,  $X_B$ . This position has a volatility of 3% per day, which is about 45% per year, or  $.03 \times \sqrt{10} = 9.49\%$  for 10 days. The corresponding VaR with a 95% confidence is

$$\text{VaR}_{.05} = \$10 \text{ million} \times .0949 \times 1.65 = \$1.565 \text{ million.}$$

Suppose the two positions are correlated with a correlation coefficient of  $\rho = .5$ . The overall daily volatility is

$$\begin{aligned}\sigma_{A+B} &= \sqrt{\sigma_A^2 + 2\rho\sigma_A\sigma_B + \sigma_B^2} \\ &= \sqrt{.02^2 + 2 \times 0.5 \times .02 \times .03 + .03^2} \\ &= 0.0436.\end{aligned}$$

This is 13.79% for 10 days, which is less than the sum of the two individual standard deviations:  $6.32\% + 9.49\% = 15.81\%$ . Again with  $h = .05$ , we find that

$$\text{VaR}_{.05}(X_A + X_B) = 10 \times .1379 \times 1.65 = \$2.275 \text{ million,}$$

which is less than the sum of the VaR's, which is \$ 2.608 million.

Note that in general

$$\text{VaR}_h(X_{A+B}) \leq \text{VaR}_h(X_A) + \text{VaR}_h(X_B) \quad (10.6)$$

(because  $\sigma_{A+B} < \sigma_A + \sigma_B$ ), demonstrating the reduction in risk of a sum compared to the sum of the risks, at least as determined by value at risk for normal random variables. Property (10.6) is termed **subadditivity**. The same conclusion holds if the  $\mu$ 's are not zero, but it does not necessarily hold for non-normal distributions. See Exercise 5. and Section 10.3.

It is easy to see that in general  $\text{VaR}_h(\frac{1}{2}X) = \frac{1}{2}\text{VaR}_h(X)$ . It follows from subadditivity that if the pension fund diversifies by holding  $\frac{1}{2}X_A + \frac{1}{2}X_B$  it would find

$$\text{VaR}_h(\frac{1}{2}X_A + \frac{1}{2}X_B) \leq \frac{1}{2}\text{VaR}_h(X_A) + \frac{1}{2}\text{VaR}_h(X_B),$$

which shows the advantage of diversification.

## Other Models

The use of the normal distribution to compute value at risk is convenient for two reasons. First, as shown earlier, there is a simple formula for the associated value at risk based only on mean and variance. Second, a position  $X$  that consists of a linear combination of separate normal components is itself normal, and hence the value at risk of the entire portfolio is again simple to calculate. However, most practitioners agree that assuming a normal distribution is not exactly realistic, especially because it drastically underestimates the severe losses represented in the extreme lower portion of the density, that is, in the lower tail. It is argued (and data verify) that the actual probabilities of extreme losses are usually substantially higher than probabilities predicted by a normal density having the same mean and variance as the true density. This fact is described by saying that the distribution has a lower **fat tail**. See Chapter 13.

This weakness can sometimes be ameliorated by assuming an alternate density, one believed to represent the tail better. However, analytic manipulation of such densities is usually cumbersome, especially when attempting to infer the overall distribution of a portfolio consisting of several complex components.

The general problem of computing VaR or similar risk measures is daunting indeed. The number of component positions can be in the thousands, many of the associated risks are nonlinear functions of market variables (such as options), and it is difficult to handle so-called “event risks” such as earthquakes or other rare events that might impact financial systems. These more complex situations are treated by analysis of scenarios, use of Monte Carlo methods, massive simulation, and extrapolation from historical records. It is a challenging but important area.

## Shortcut for Discrete Distributions

Sometimes a position  $X$  is defined by a finite set of points, each with its own probability and  $X$ -value. The individual values may be defined directly, or often they are the final nodes in a tree or lattice model. In any case, there is a simple rule for determining the associated value at risk. One begins by ordering the (final) nodes according to their  $X$ -value, from lowest to highest. Then, a running sum of total probability is formed,

starting with the lowest  $X$ -value node. The critical point is where the sum first equals or exceeds the designated fraction  $h$ . If that sum is strictly greater than  $h$ , the  $X$ -value at that node (with sign reversed) is the value at risk. If the sum at the critical node is exactly equal to  $h$ , then one must continue until the sum is strictly greater than  $h$  (at the next node with positive probability). The value at that node (with sign reversed) is then the value at risk. For example, suppose the two lowest values of  $X$  are  $-\$100$  and  $-\$70$  and that each has probability 10%. Using  $h = .10$ , the VaR would be  $\$70$ ; while using  $h = .09$ , the VaR would be  $\$100$ . This shortcut method is helpful in studying examples with discrete distributions that illustrate various properties of value at risk.

## Empirical Approach for Market Risk\*

The probability distribution of a position  $X$  can often be inferred directly from historical price data of the assets in the position. For estimates of daily returns, the basic empirical method is quite straightforward.

One begins with a database of daily rates of return for each asset  $i$  in the portfolio over a long period, say, 500 days. These returns are indexed  $r_{ik}$ , where  $k = 1, 2, \dots, 500$  is a day index. Suppose that for each asset  $i$ ,  $x_i$  is the current dollar amount held in asset  $i$ . Then  $\sum_{i=1}^n x_i r_{ik}$  is the amount of value change that the portfolio will have tomorrow if today's returns happen to be the same as they were on day  $k$ . By making this calculation for each day  $k$ , we obtain 500 hypothetical portfolio change values  $\sum_{i=1}^n x_i r_{ik}$  for tomorrow. A probability of 1/500 can be assigned to each of these possible occurrences, and the results could be laid out as a histogram, displaying the density for tomorrow's value. If the confidence level is set at 98%, the return that is 2% from the bottom (that is the 10th lowest) is critical. In view of the shortcut method, we should use the 11th-lowest rate of return. The hypothetical loss of today's current value when subjected to this 11th-worst rate of return is the value at risk for 1 day at confidence level of 98%.

**Example 10.4 (Ten days)** Suppose you have recorded the returns of every asset in your portfolio on each of 10 consecutive days. From the asset portfolio weights, you compute the corresponding hypothetical daily returns, which are shown across the first line of Table 10.1, day by day. You then sort these returns by their magnitude, from lowest to highest, and enter them in the second row of the table. Suppose you wish to find the value at risk for  $h = .20$ . Each element in the row accounts for 10% of the total probability. Summing these probabilities starting at the left, the total probability of exactly .20 is attained at the second entry. According to the shortcut method, you must go one more step, which is  $-.07$ . For an initial value of \$100,000,

**TABLE 10.1**  
**TEN DAYS OF RETURN DATA**

Unsorted	.05	-.01	-.13	.06	.08	-.15	.09	-.04	-.07	.03
Sorted	-.15	-.13	-.07	-.04	-.01	.03	.05	.06	.08	.09

the value at risk is \$7,000. Thus, there is less than or equal to a 20% chance that you will lose more than \$7,000.

Notice that even with 500 samples, it is only the rankings among a small group of low returns that determine the value at risk. It seems intuitively clear that this small sample is not likely to be terribly representative. However, the basic method described can be improved by supplementing it with sophisticated methods for inferring values in the tail. These methods include testing whether certain standard distributions fit the data well and using a good fit to smooth the result.

In any event, the method must be augmented if it is desired to determine the value at risk for a period longer than a single day. One way to extend the value to  $N$  days is simply to multiply the VaR by  $\sqrt{N}$ , which is exact in the case of normal distributions with means of zero. Alternatively, one could construct, from the original historical data, blocks of length  $N$ , recording the  $N$ -day returns. However, if nonoverlapping blocks are used, they provide only  $500/N$  independent points. For  $N = 5$ , for example, this provides only 100 points, and the 98% confidence level would be the second (or third) smallest of these returns—clearly an imprecise measure.

An alternative is to use overlapping blocks of data, the first being for the period  $k = 1, 2, \dots, 10$ , the second being for the period  $k = 2, 3, \dots, 11$ , and so forth. There are close to 500 of these, but they are not independent and hence are not much better than using nonoverlapping blocks.

## 10.3 Criticisms of VaR

Value at risk has had wide acceptance in the financial world; in fact, since the Basel Accords of 1988, it has provided an accepted standard for bank regulation. Initially it was applied to credit risk, but later, in 1996, the Basel requirements were applied to market risk as well, using a confidence level of 99% and a 10-day horizon. The Basel Accords allow some flexibility in how value at risk is determined, and there are indeed many methods in use, both analytic and empirical. In 2010 the Accords were updated again, requiring banks to hold greater capital reserves.

The broad acceptance of value at risk as a universal standard raises it to a very high level of importance in the theory of risk measurement and control. Nevertheless, despite the ubiquitous nature of value at risk, it has been criticized as having several troublesome shortcomings. Three of the most important of these are outlined here.

### Diversification Failure

Example 10.2 showed that when distributions are jointly normal, the value at risk of a sum is less than the sum of the individual value at risk figures, a property termed *subadditivity*. That is, for a fixed loss tolerance  $h$ ,

$$\text{VaR}_h(X_1 + X_2) \leq \text{VaR}_h(X_1) + \text{VaR}_h(X_2). \quad (10.7)$$

It is also true for the case of normally distributed positions, that the value at risk of a weighted average of positions is less than or equal to the corresponding weighted average of the individual value at risk figures, reflecting the advantage to diversification. However, this diversification advantage is not necessarily reflected in VaR for distributions that are not normal, as illustrated by the following example.

**Example 10.5 (Diversification failure)** Consider the position

$$X_1 = \begin{cases} 1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2 \end{cases}. \quad (10.8)$$

Set  $h = .50$ . To compute  $\text{VaR}_{.50}(X_1)$ , we use the shortcut method for discrete probabilities. Starting at the value  $-1$ , we have exactly  $.50$ , so VaR is determined at the next higher node, which is  $1$ . Changing the sign gives  $\text{Var}_{.50}(X_1) = -1$ .

Now suppose that  $X_2$  has exactly the same structure as  $X_1$  but is independent of  $X_1$ . Consider the diversified position  $X$  made up of equal parts of  $X_1$  and  $X_2$ . This has payoff

$$X = \frac{X_1 + X_2}{2} = \begin{cases} 1 & \text{with probability } 1/4 \\ 0 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/4 \end{cases}. \quad (10.9)$$

Again, using the shortcut method, we start at the lowest-valued point and move upward accumulating probability. The value  $.50$  is exceeded at  $0$  with an accumulated probability of  $3/4$ . Since this is greater than  $.50$ , it defines VaR as  $\text{VaR}_{.50}(X) = 0$ . Thus for  $h = 0.50$  we find that

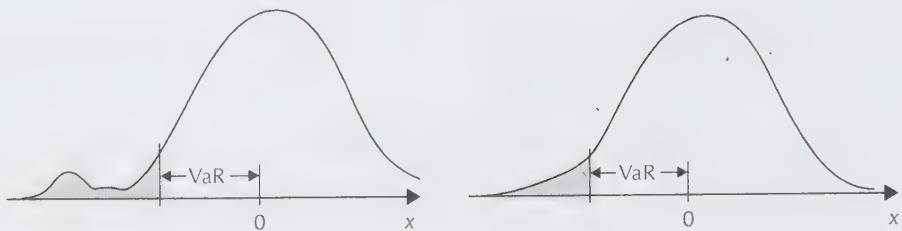
$$\text{VaR}_h\left(\frac{X_1 + X_2}{2}\right) > \frac{1}{2}\text{VaR}_h(X_1) + \frac{1}{2}\text{VaR}_h(X_2),$$

which implies that the diversified position is riskier than the weighted sum of each of its component VaR's.

The example shows that although diversified portfolios have advantages over their nondiversified counterparts (as, for example, having lower variance), the corresponding VaR figures may not reflect these advantages. Although this phenomenon does not always occur, the fact that it can occur is considered by many researchers to be a shortcoming of value at risk as a risk measure. The phenomenon is related to another, which is illustrated next.

## Poor Assessment of Risk

Perhaps the most severe criticism of value at risk is that it does not account for the fact that some losses may be much greater than the value at risk. Value at risk tells us the minimum loss in the  $h\%$  worst cases; it does not say anything about how bad losses might be in the  $h\%$  worst cases. Value at risk gives only limited information about the low end of the density.



**FIGURE 10.4** Value at risk does not distinguish between patterns of losses in the lower tail.

**Example 10.6 (Dangerous situation)** Suppose  $h = .05$ . Consider the two positions

$$X_1 = \begin{cases} 100 & \text{with probability 92\%} \\ -10 & \text{with probability 6\%} \\ -15 & \text{with probability 2\%} \end{cases}$$

$$X_2 = \begin{cases} 100 & \text{with probability 92\%} \\ -10 & \text{with probability 6\%} \\ -100 & \text{with probability 2\%} \end{cases}.$$

Both positions have  $\text{VaR}_{.05} = 10$ . However, intuitively, there is more risk associated with  $X_2$  than with  $X_1$ . Value at risk does not distinguish between these two positions.

This effect is illustrated in Figure 10.4, which shows two continuous probability densities, both of which have the same value at risk. However, the situation on the left seems riskier than that on the right because there are heavy losses on the left. In general, value at risk shows the upper limit point of the fraction  $h$  of worst cases but says nothing about how bad the losses worse than that can be.

## Discontinuous Value

Another criticism of value at risk is that it can be very sensitive to the choice of the risk tolerance  $h$ , especially for discrete-valued  $X$ 's.

**Example 10.7 (Jumpy VaR)** Consider  $X$  defined as

$$X = \begin{cases} 1 & \text{with probability 1/2} \\ -1 & \text{with probability 1/2} \end{cases}.$$

The VaR at  $h = .50$  is  $-1$ . However, the VaR at  $h = .49$  is  $1$ .

The fact that  $\text{VaR}_h(X)$  may be discontinuous with respect to  $h$  means that the risk assessment may change significantly if the required risk tolerance is changed only slightly.

## 10.4 Coherent Risk Measures

Discussion and evaluation of risk-measure properties can be put on an objective framework by spelling out exactly what properties a risk measure should possess ideally. This idea was first carried out in a paper by Artzner, Delbaen, Eber, and Heath [6], and we shall outline their framework here.

Consider a class of random variables (defined on a common set of possible outcomes), with each random variable representing, as before, the profit  $X$  of some position at a fixed time. It is possible to form linear combinations of positions in this class, and there is a risk-free asset. A **risk measure** for this class is a function  $\rho$  mapping  $X$ 's into real numbers. A position  $X$  is considered **acceptable** according to the risk measure  $\rho$  if  $\rho(X) \leq 0$ . Hence, a risk measure partitions the possible positions into two sets: those that are acceptable and those that are not.

A risk measure is said to be **coherent** if it satisfies the following four axioms.<sup>6</sup>

1. **Translation Invariance (T)**. For all  $X$  and all risk-free capital payoffs  $C$ , there holds  $\rho(X + C) = \rho(X) - C$ .
2. **Subadditivity (S)**. For all  $X_1, X_2$  there holds  $\rho(X_1 + X_2) \leq \rho(X_1) + \rho(X_2)$ .
3. **Positive Homogeneity (PH)**. For all  $\lambda \geq 0$  and all  $X$ , there holds  $\rho(\lambda X) = \lambda \rho(X)$ .
4. **Monotonicity (M)**. For all  $X_1, X_2$  with  $X_1 \leq X_2$ , there holds  $\rho(X_2) \leq \rho(X_1)$ .

A simple example of a coherent risk measure is  $\rho(X) = E[-X]$ , where  $E$  is the ordinary expected value. In fact, it is perhaps a measure that one might seriously consider first. More complex coherent risk measures can be constructed by combining various expected-value measures using different probabilities, such as those associated with pessimistic random forecasts, as discussed in Section 10.6.

In general, the first axiom, translation invariance, implies that the risk measure is based on units of capital. It implies in particular that if a position  $X$  has level of risk  $\rho(X) > 0$ , so that it is not acceptable, it can be made acceptable by adding a risk-free asset  $C$  with payoff  $\rho(C)$ . Mathematically,  $\rho(X + \rho(X)) = 0$ .

The next axiom, subadditivity, ensures for one thing that if each of two positions is acceptable, then their sum is also acceptable. In general, together with axiom 1, it implies that the sum of two positions can be made acceptable by the addition of cash equal to or less than the sum of the cash required to make each of the components acceptable.

Axiom 3, positive homogeneity, states that risk scales up proportionally with the size of the position. This axiom, together with subadditivity, gives

$$\rho(\alpha X_1 + (1 - \alpha) X_2) \leq \alpha \rho(X_1) + (1 - \alpha) \rho(X_2)$$

for all  $\alpha$  with  $0 \leq \alpha \leq 1$ , which mathematically means that  $\rho$  is a convex function and financially means that diversification is never disadvantageous.

---

<sup>6</sup> It is assumed that  $\rho(X) \neq \infty$ .

Finally, monotonicity simply states that if a first position has a larger payoff than a second in every possible situation, then the first is considered less risky than the second.

It will be noticed that value at risk satisfies three of the four axioms but fails to satisfy subadditivity. We immediately conclude that value at risk is *not* a coherent risk measure.

You should be able to verify that the simple measure  $\rho(X) = E[-X]$  satisfies all four axioms.

## 10.5 Conditional Value at Risk

The publication of the axiomatic framework for assessing risk measures initiated a search for alternative measures in the spirit of value at risk, in that they consider downside risk but are coherent. In particular, there was great interest in finding an alternative to VaR that would account for the pattern of losses that are more severe than the specific loss defined by value at risk.

A measure that satisfies these conditions at least partially is the expected value of loss conditional on the loss's being at least equal to the VaR value. In other words, rather than asking what the minimum loss in the worst  $h$ -fraction of cases is, we ask what the expected loss in the worst  $h$ -fraction of cases is. This measure is termed the **conditional value at risk** (CVaR). In terms of probability, this new measure is essentially the expected value of the area of the tail of the density up to the  $h$ -quantile. Specifically, for distributions that are strictly increasing without jumps, CVaR is defined by<sup>7</sup>

$$\text{CVaR}_h(X) = E\left\{-X \mid X \leq -\text{VaR}_h(X)\right\}. \quad (10.10)$$

Conditional value at risk satisfies the inequality

$$\text{CVaR}_h(X) \geq \text{VaR}_h(X)$$

because the losses in the  $h$ -tail are at least as large as  $\text{VaR}_h(X)$ , and  $\text{CVaR}_h(X)$  is the average of these. Measure (10.10) has alternatively been referred to as **tail value at risk** (TVaR), **expected shortfall**, and **conditional tail expectation** (CTE). An alternative measure is the **average value at risk**, defined as

$$\text{AVaR}_h(X) = \frac{1}{h} \int_0^h \text{VaR}_u(X) du. \quad (10.11)$$

However, as discussed in Exercise 8, AVaR is equivalent<sup>8</sup> to CVaR.

The definition of CVaR given by measure (10.10) is coherent when the probability distribution is strictly increasing and without jumps. However, to ensure coherence for general distributions, one must more carefully define the range over which to compute the expected value. In particular, adjustments must be made if the value at risk

<sup>7</sup> That is, the expected value of  $-X$ , given that  $X \leq -\text{VaR}$ .

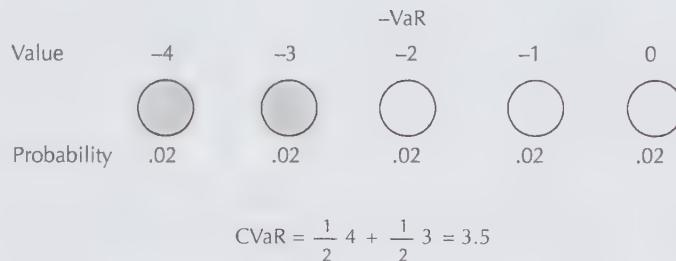
<sup>8</sup> Valid when the expected value in measure (10.10) and the integral in measure (10.11) are defined properly.

point is a probabilistic **atom** (that is, if it itself has positive probability), for such points cause jumps in the distribution. The issue is that the (conditional) expected value should be computed over the tail region having probability exactly equal to  $h$ . For distributions with atoms, exact equality may not occur. In those cases, an adjustment is made by splitting the atoms appropriately. This is easy to understand by study of the following examples.

**Example 10.8 (No atom splitting needed)** Suppose the leftmost tail of  $X$  is discrete-valued with the following values:

$$X = -4 \ -3 \ -2 \ -1 \ 0 \ 1 \dots$$

All of these have probability of 2%. We wish to find  $\text{CVaR}_h(X)$  with  $h = .04$ . The situation is shown in Figure 10.5.



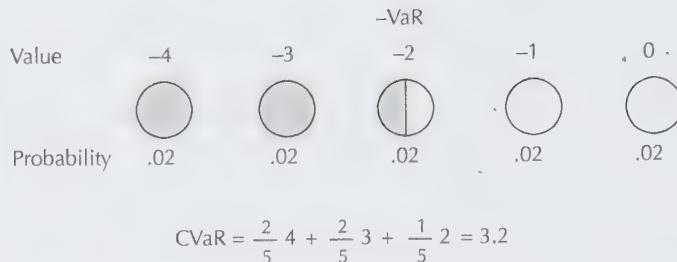
**FIGURE 10.5** With  $h = .04$ , we find  $\text{VaR} = 2$ . The two atoms below  $-\text{VaR}$  have a total probability of .04, so they define the region for computing the conditional expectation. Since they both have the same probability, they are weighted equally in the conditional-expectation calculation.

We immediately see that  $\text{VaR}_{.04}(X) = 2$ , because the total probability of the two shaded atoms is exactly .04 (and hence we must proceed to the next node to get  $\text{VaR}$ ). Next, we find the expected value of  $-X$  over the shaded nodes, which together have a probability of .04. Since these both have the same probability, they are weighted equally in the conditional expectation. Hence  $\text{CVaR}_{.04}(X) = \frac{1}{2} 4 + \frac{1}{2} 3 = 3.5$ .

**Example 10.9 (Splitting the atom)** Suppose for the same  $X$  as in Example 10.8 that we wish to use  $h = .05$ . This situation is shown in Figure 10.6.

The two atoms below  $-\text{VaR}$  have a total probability of 4%, but  $h = 5\%$ , so we must split the atom at the  $\text{VaR}$  point and assign 1% to the lower tail to get the full 5%. This defines the region for computing the conditional expectation. The relative weights of the atoms to be used in the conditional expectation are now  $2/5, 2/5, 1/5$ , and thus  $\text{CVaR}_{.05}(X) = \frac{2}{5} 4 + \frac{2}{5} 3 + \frac{1}{5} 2 = 3.2$ .

It is not difficult to verify that  $\text{CVaR}$  satisfies most of the axioms required of coherency. The one that is difficult to establish is subadditivity, but this has in fact been



**FIGURE 10.6** With  $h = .05$ , we again find  $\text{VaR} = 2$ . However, the atom at the  $\text{VaR}$  point must be split so that the probabilities of the shaded atoms and partial atom sum to .05.

proved. So CVaR (when splitting is used if necessary) is coherent even for discrete probability distributions. It also overcomes many of the objections to ordinary value at risk. It accounts for the shape of losses that are more severe than VaR, and it is continuous with respect to the confidence level. Since CVaR is more conservative than VaR, it requires that more capital be added to a risky position to make it acceptable than would be required by VaR. An attractive feature is that it is convex with respect to  $X$ . This means that if one is seeking, for example, to design a portfolio protected somewhat from loss while satisfying various constraints on portfolio weights, the formal objective of minimizing CVaR subject to those constraints will have convexity properties that facilitate optimization.

## 10.6 Coherent Characterization\*

The original paper defining coherent risk measures also presented a universal characterization of such measures. Although the characterization is rarely used in practice to design or implement measures, it provides valuable insight into what it means for a risk measure to be coherent.

In what follows, we make explicit reference to the set  $\Omega$ . This set is assumed to be finite, containing, say,  $n$  elements that can be thought of as possible (financial) states of the world. A set of “real” probabilities is defined on  $\Omega$ . A situation  $X$  is a random variable that takes values on the elements of  $\Omega$ . The real probabilities associated with the states remain fixed, but for the characterization in this section we consider other possible probability assignments. Such an assignment is a vector of  $n$  nonnegative numbers that sum to 1. We denote such an assignment vector by  $p$ .

Once  $p$  is specified, the corresponding expected value is defined as  $E_p(X) = \sum_{i=1}^n p_i X_i$ , where  $p_i$  and  $X_i$  are the components of  $p$  and  $X$ . In other words, expected value is computed using the  $p_i$ ’s as the probabilities.

For the following characterization we restrict the allowable  $p$ ’s to a **family**  $\mathcal{P}$ . Such a  $\mathcal{P}$  is a subset of  $n$ -dimensional space  $R^n$ . Each family defines a risk measure, as explained in the theorem.

**Theorem 10.1 Coherent Characterization** Suppose  $\Omega$  is a finite set. A risk measure  $\rho$  is coherent if and only if there exists a family  $\mathcal{P}$  of probabilities on  $\Omega$  such that<sup>9</sup>

$$\rho(X) = \max_{p \in \mathcal{P}} \left\{ E_p[-X] \right\}. \quad (10.12)$$

We can quickly relate several simple risk measures to the characterization defined by this theorem.

1. The ordinary expected value  $\rho(X) = E[-X]$  is the measure where the family  $\mathcal{P}$  consists of just one vector, the actual (“true”) probability vector on the states of the world.
2. The worst-case measure  $\rho(X) = -\min(X)$  is obtained by taking the family  $\mathcal{P}$  to be the entire set of all possible probabilities defined on  $\Omega$ . Then, if the worst case is on state  $i$ , the vector  $p$  chosen by the max in equation (10.12) will have all components zero except the  $i$ th, which will be 1.
3. The forecasting measure  $\rho(X)$  is defined by including only a single probability vector in  $\mathcal{P}$ : It may be a certain expert’s best guess of the “true” probabilities.
4. The stress test measure is defined by a family  $\mathcal{P}$  consisting of one or more pessimistic probability vectors defining various possible futures.

It is easy to prove the “if” portion of the characterization theorem; that is, it is easy to show that the measure defined by Theorem 10.1 is coherent. Clearly it is translation invariant, because adding a constant to  $X$  will simply subtract that constant from the resulting risk measure. Likewise, monotonicity and positive homogeneity are clear. Subadditivity is the most difficult, but even that is not hard. For any  $X_1$  and  $X_2$  we have

$$\max_{\mathcal{P}} \left\{ E_p[-X_1 - X_2] \right\} \leq \max_{\mathcal{P}} \left\{ E_p[-X_1] \right\} + \max_{\mathcal{P}} \left\{ E_p[-X_2] \right\}$$

because the maximum of a sum is less than the sum of the maxima.

**Example 10.10 (CVaR characterization)** Let us see how in a specific case the coherent measure CVaR can be represented as in the theorem. Consider a situation where there are five states of the world, each with a probability of .2. We wish to determine the CVaR for any  $X$  defined on the five states. For example, suppose we wish to compute  $\rho(X) = \text{CVaR}_{40\%}(X)$  for various  $X$ ’s. As a particular case, take

$$X_1 = \begin{cases} 10 \\ 5 \\ 3 \\ -4 \\ -9 \end{cases}.$$

The average loss of the 40% worst cases is  $\text{CVaR}_{40\%}(X_1)$ . (With this choice of  $h$  and the given probabilities, there will be no need to split atoms.) For the given  $X_1$ ,  $\text{VaR} = -3$ ,

<sup>9</sup> Again “sup” may be more appropriate than “max” here.

**TABLE 10.2**  
**PROBABILITY FAMILY FOR EXAMPLE 10.10**

State	Probability									
1	.5	.5	.5	.5	0	0	0	0	0	0
2	.5	0	0	0	.5	.5	.5	0	0	0
3	0	.5	0	0	.5	0	0	.5	.5	0
4	0	0	.5	0	0	.5	0	.5	0	.5
5	0	0	0	.5	0	0	.5	0	.5	.5

Each column represents a set of probabilities for the five states.

and CVaR is the average of the loss from states 4 and 5; namely,  $\text{CVaR}_{40\%}(X_1) = [0.5 \times 4 + 0.5 \times 9] = 6.5$ .

According to the theorem, there is a family of probability measures that define  $\text{CVaR}_{40\%}(X)$  in general for this world  $\Omega$ , that is, for any  $X$ , not just  $X_1$ . In this example, that family consists of all those probability vectors with .5 probability in two states and zeros in the rest. The entire family is shown in Table 10.2.

For the foregoing  $X_1$ , the last probability set in the table is the one that yields the largest expected value of  $-X_1$ , and hence the risk measure is just the expected loss using this vector; again 6.5.

Now consider the payoff

$$X_2 = \begin{cases} -8 \\ 12 \\ 7 \\ 15 \\ 4 \end{cases}$$

In this case the probability vector that produces the largest expected loss is the fourth one in the table, the one with .5 probability for states 1 and 5. Hence

$$\text{CVaR}_{40\%}(X_2) = [.5 \times 8 - .5 \times 4] = 2.$$

## 10.7 Convexity\*

The appropriateness of the axiom of positive homogeneity required for coherence has often come into question. It is argued that increasing a position by a large positive factor should lead to a risk level higher than simple proportionality would dictate. This has motivated some theorists to drop both this axiom and the subadditivity axiom and to replace them with the single convexity axiom, as stated next:

(C) A risk measure  $\rho$  is **convex** if for all  $X_1$  and  $X_2$

$$\rho(\alpha X_1 + (1 - \alpha) X_2) \leq \alpha \rho(X_1) + (1 - \alpha) \rho(X_2)$$

for all  $\alpha$  with  $0 \leq \alpha \leq 1$ .

In the original set, the axioms of positive homogeneity and subadditivity together imply convexity. However, convexity does not imply positive homogeneity; hence, replacing positive homogeneity by convexity loosens the axiom system somewhat. A similar but slightly more complex characterization theorem is available for this new system. Furthermore, the presence of convexity facilitates optimization and opens the possibility of interesting duality results.

## 10.8 Summary

In some situations, such as monitoring a bank's financial position, it is reasonable to focus almost exclusively on the risk of loss rather than on the possibility of gain. In mathematical terms, this means focusing on the lower tail of the probability density of the position's payoff  $X$ . A popular measure is **value at risk** (VaR), which, at a given risk tolerance  $h$ , is the smallest number  $V$  such that the probability of a loss greater than  $V$  is no more than  $h$ .

A basic property of value at risk is that, if a risk-free payoff of amount  $C$  is added to the final value of the position, the value at risk will be reduced by exactly  $C$ . Since the Basel Accords of 1998, value at risk has formed an international standard for assessing the soundness of financial institutions.

Implementation of value at risk requires estimation of the lower tail of the payoff density. This is reasonably feasible when the risk of loss is short-term market risk over a few days. Determination of VaR can be much more challenging for nonmarket risk, such as credit risk. In the case of short-term market risk, if the returns are normally distributed, then there is a simple analytic formula for computing VaR. In most cases, however, extensive statistical analysis must be conducted to estimate the critical features of the loss tail.

Despite its nearly universal acceptance, value at risk has been criticized as not being fully representative of potential loss. It provides only the simplest information about the tail of the distribution, failing to distinguish heavy losses from light losses beyond VaR; it may be discontinuous with respect to the risk tolerance if there are jumps or flat spots in the distribution; and it does not reflect an advantage to diversification.

A risk measure is **coherent** if it satisfies a set of four specific axioms. Value at risk satisfies only three of these, failing to satisfy subadditivity. An alternative coherent risk measure commonly termed **conditional value at risk** (CVaR) is defined as the expected value of loss conditional on the loss being below the lower  $h$ -quantile tail, where  $h$  is the risk tolerance and  $1 - h$  is the confidence level. This measure has several advantages compared to value at risk.

## Exercises

1. (A big position) The position  $X$  has a 1-year probability density function that is normal with a mean of \$100 million and a standard deviation of \$50 million. What is the value at risk at the 99% confidence level?

2. (Normal examples) Suppose  $X$  is a normal with zero mean and standard deviation of \$10 million.

(a) Find the value at risk for  $X$  for the risk tolerances  $h = 0.01, 0.02, 0.05, 0.10, 0.50, 0.60$ , and 0.95.

(b) Is there a relation between VaR for values of  $h \leq 0.50$  and values for  $h \geq 0.50$ ?

3. (Uniform) Consider the position  $X$  that has a uniform probability density between  $-40$  and  $60$ . Find  $\text{VaR}_h(X)$  for all  $h$ ,  $0 \leq h \leq 1$ .

4. (Equivalences\*) General equivalent formulas for VaR are shown here. Assuming the first one, argue that the next two are equivalent to it. [Hint: Consider a distribution  $X$  that at a certain point jumps up. Then examine points to the right of the jump for a formula such as (a) or to the left for one such as (c).]

$$(a) \text{VaR} = -\inf \left\{ m : P[X \leq m] > h \right\}$$

$$(b) \text{VaR} = \sup \left\{ y : P[X \leq -y] > h \right\}.$$

$$(c) \text{VaR} = -\sup \left\{ m : P[X < m] \leq h \right\}$$

$$(d) \text{VaR} = \inf \left\{ y : P[X < -y] \leq h \right\}.$$

$$(e) \text{VaR} = \inf \left\{ -m : P[X < m] \leq h \right\}.$$

$$(f) \text{VaR} = \inf \left\{ C : P[X + C < 0] \leq h \right\}.$$

5. (Subadditivity) Suppose  $X_1$  and  $X_2$  are jointly normal positions with parameters  $\mu_1, \mu_2, \sigma_1, \sigma_2, \rho_{12}$ . Show that

$$\text{VaR}_h(X_1 + X_2) \leq \text{VaR}_h(X_1) + \text{VaR}_h(X_2).$$

6. (AVaR calculation) Find  $\text{AVaR}_h(X)$  for the  $X$  of Exercise 3. [Also see Exercise 9.]

7. (A standard measure) Let  $a > 0$  and consider the risk measure  $\rho(X) = E[-X] + a\sigma(X)$ . Is this a coherent risk measure? If not, which axioms are violated?

8. (CVaR–AVaR equivalence) Let  $X$  be a position with a probability distribution  $F$  that is strictly increasing and smooth. Let  $f(x) = F'(x)$  be the associated probability density.

(a) Verify that

$$\text{CVaR}_h(X) = -\frac{1}{h} \int_{-\infty}^{-\text{VaR}_h(X)} xf(x)dx \quad (10.13)$$

(b) For any  $u \in (0, 1)$  let  $x = F^{-1}(u)$  be the value of  $X$  that defines the  $u$ -quantile of  $X$ . Conversely, for any specific value  $x$  of  $X$ , we have  $u = F(x)$  as the quantile value associated with  $x$ . Using the change of variable  $u = F(x)$  in equation (10.13), show that

$$\text{CVaR}_h(X) = -\frac{1}{h} \int_0^h F^{-1}(u)du. \quad (10.14)$$

(c) Interpret the right-hand side of equation (10.14) to obtain

$$\text{AVaR}_h(X) = -\frac{1}{h} \int_0^h F^{-1}(u)du. \quad (10.15)$$

and hence conclude that  $\text{CVaR}_h(X) = \text{AVaR}_h(X)$ . [Equation (10.15) can be extended to be coherent in general (e.g., with jumps in  $F$ ) because the upper limit of the integral will force atoms to be split if necessary so that the integration is strictly over the range zero to  $h$ .]

9. (CVaR uniform) Find  $\text{CVaR}_h(X)$  for the linear case of Exercise 3. [Also see Exercise 6.]
10. (Diversification advantage) For Example 10.5 with  $h = 50\%$ , show explicitly that CVaR rewards diversification.
11. (Modified confidence) Modify Example 10.10 by changing the confidence level to 70% (that is, with  $h = 30\%$ ).
  - (a) Find the conditional value at risk for  $X_1$  of the example.
  - (b) Describe the family  $\mathcal{P}$  of probability distributions that characterize this risk measure.
12. (Change  $h^*$ ) Suppose that a bank has position  $X$  that has a normal probability density. The value at risk is known to be  $V$  at the loss tolerance  $h$ . The bank plans to take some capital out of reserve, which would increase the VaR. The bank argues that the change would be acceptable if the required confidence level were only slightly increased.
  - (a) Develop a formula that expresses implicitly, in terms of the inverse standard normal distribution function  $F_N^{-1}$ , the amount that the loss tolerance must change to compensate for a change  $\Delta$  in capital. [Hint: Express the result as the difference of two quantities.]
  - (b) Approximate the difference in (a) in terms of a derivative, and use the fact that  $\frac{d}{dh}F_N^{-1}(h) = 1/f(x)$ , where  $x = F_N^{-1}(h)$  and where  $f$  is the probability density of a standardized normal random variable.
  - (c) Using the fact that  $h \approx 0$ , find the amount of the required change as a fraction of  $\sigma$ .
  - (d) For  $h = 1\%$  and  $\bar{X} = -100$ ,  $\sigma = 20$ , find the required new value of  $h$  for an addition of  $-10$  to  $\bar{X}$ .
  - (e) Verify the result of part (d) explicitly (without the approximation).

## References

For general background on value at risk at an introductory level, see [1]. For a detailed account see [2]. A careful overview of a more advanced nature is [3]. The standard reference on the statistics of tail events is [4]. An advanced treatment of risk measures is presented in [5].

The concept of a coherent risk measure defined by axioms was put forth by Artzner, Delbaen, Eber, and Heath in [6] and [7]. The theory leading to various representations of conditional value at risk and proof of their coherence are found in [8], [9], [10]. For good overviews of CVaR, see [12], [13], and [14]. The study of convex risk measures and associated duality theory is contained in [11].

An excellent advanced overview of CVaR and its advantage in the formulation and solution of optimization problems is by Rockafellar [12].

1. Jorion, Philippe (1997), *Value at Risk*, Irwin, Chicago.
2. Duffie, D., and J. Pan (Spring 1997), "An Overview of Value at Risk," *Journal of Derivatives*, 4, no. 3, 7–49.
3. Giesecke, Kay (2009), Class notes, Department of Management Science & Engineering, Stanford University, Stanford, CA.
4. Embrechts, Paul, Claudio Klüppelberg, and Thomas Mikosch (2001), *Modelling Extremal Events for Insurance and Finance*, Springer Verlag, Berlin, corrected third printing.

5. Föllmer, Hans, and Alexander Schied (2004), *Stochastic Finance: An Introduction in Discrete Time*, 2nd ed., de Gruyter, Berlin.
6. Artzner, Ph., F. Delbaen, J.-M. Eber, and D. Heath (November 1997), "Thinking Coherently," *RISK*, **10**, 69–71.
7. Artzner, Ph., F. Delbaen, J.-M. Eber, and D. Heath (1999), "Coherent Measures of Risk," *Mathematical Finance*, **9**, 203–228.
8. Acerbi, C., C. Nordiko, and C. Sirtori (2001), "Expected Shortfall as a Tool for Financial Risk Management," Working paper. Download at [www.gloriamundi.org](http://www.gloriamundi.org).
9. Acerbi, C. and D. Tasche (2002), "On the Coherence of Expected Shortfall," *Journal of Banking and Finance*, **26**, 1487–1503.
10. Pflug, G. (2000), "Some Results on Value-at-Risk and Conditional-Value-at-Risk," in S. Uryasev, ed., *Probabilistic Constrained Optimization: Methodology and Applications*, Kluwer Academic, Norwell, MA.
11. Föllmer, H., and A. Schied (2002), "Convex Measures of Risk and Trading Constraints," *Finance and Stochastics*, **6**, no. 4, 429–447.
12. Rockafellar, R. Tyrrell (2007), "Coherent Approaches to Risk in Optimization under Uncertainty," *Tutorials in Operations Research*, INFORMS, 38–61.
13. Uryasev, S., S. Sarykalin, and G. Serraino (2008), "VaR vs. CVaR in Risk Management and Optimization: Methodology and Applications," INFORMS Tutorial.
14. Rockafellar, R. Tyrrell, and S. Uryasev (2002), "Conditional-Value-at-Risk for General Loss Distributions," *Journal of Banking and Finance*, **26**, 1443–1471.

# 11

## GENERAL PRINCIPLES

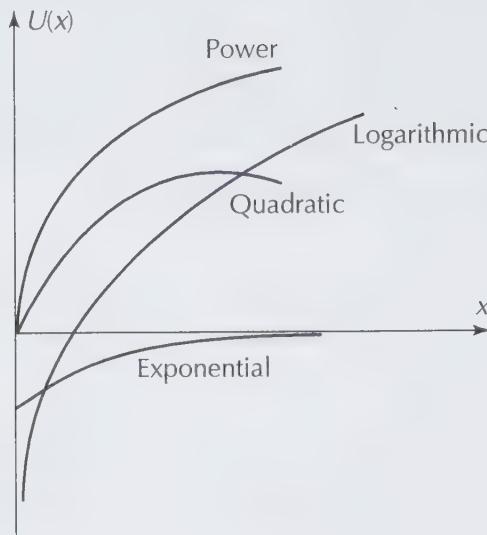
### 11.1 Introduction

Fundamentally, there are two ways to evaluate a random cash flow: (1) directly, using measures such as expected value and variance; and (2) indirectly, by reducing the flow to a combination of other flows which already have been evaluated. This chapter focuses on these two approaches, showing how they apply to single-period investment problems—and showing how they work together to produce strong and useful pricing relationships.

This chapter is more abstract than the previous chapters and serves primarily as preparation for the study of general multiperiod problems in Parts 3 and 4. The reader may wish to skip ahead to Chapter 12 (or even Chapter 13) since most of the material in Part 3 can be understood without studying this chapter. One strategy is to study the first part of this chapter—the first five sections, which cover expected utility theory. Then later, when approaching Part 4, the reader can come back to the second part of this chapter to study general pricing theory. Other readers may wish to study this chapter in sequence, for it is a logical culmination of the single-period framework.

### 11.2 Utility Functions

Suppose that, sitting here today, you have a number of different investment opportunities that could influence your wealth at the end of the year. Once you decide how



**FIGURE 11.1 Some popular utility functions.** Utility functions should increase with wealth, since greater wealth is preferred to less wealth. Functions with simple analytic forms are convenient for representation and analysis.

to allocate your money among the alternatives, your future wealth is governed by corresponding random variables. If the outcomes from all alternatives were certain, it would be easy to rank the choices—you would select the one that produced the greatest wealth. In the general random case, however, the choice is not so obvious. You need a procedure for ranking random wealth levels. A utility function provides such a procedure.

Formally, a **utility function** is a function  $U$  defined on the real numbers (representing possible wealth levels) and giving a real value. Once a utility function is defined, all alternative random wealth levels are ranked by evaluating their expected utility values. Specifically, you compare two outcome random wealth variables  $x$  and  $y$  by comparing the corresponding values  $E[U(x)]$  and  $E[U(y)]$ ; the larger value is preferred.

The specific utility function used varies among individuals, depending on their individual risk tolerance and their individual financial environment. The simplest utility function is the linear one  $U(x) = x$ . An individual using this utility function ranks random wealth levels by their expected values. This utility function (and an individual who employs it) is said to be **risk neutral** since, as will become clear later, no account for risk is made. Other utility functions do account for risk.

The one fairly general restriction is that the utility function be an increasing, or strictly increasing, continuous function. [ $U$  is **increasing** if for real numbers  $x$  and  $y$  with  $x \geq y$  it follows that  $U(x) \geq U(y)$ .  $U$  is **strictly increasing** if  $x > y$  implies  $U(x) > U(y)$ .] Other than a restriction of this sort, the utility function can, at least in

theory, take any form. In practice, however, certain standard types are popular. Here are some of the most commonly used utility functions (see Figure 11.1):

### 1. Exponential

$$U(x) = -e^{-ax}$$

for some parameter  $a > 0$ . Note that this utility has negative values. This negativity does not matter, since only the *relative* values are important.

### 2. Logarithmic

$$U(x) = \ln(x).$$

Note that this function is defined only for  $x > 0$ . In fact, if there is any positive probability of obtaining an outcome of 0, the expected utility will be  $-\infty$ .

### 3. Power

$$U(x) = bx^b$$

for some parameter  $b \leq 1, b \neq 0$ . This family includes (for  $b = 1$ ) the risk-neutral utility.

### 4. Quadratic

$$U(x) = x - bx^2$$

for some parameter  $b > 0$ . Note that this function is increasing only for  $x < 1/(2b)$ .

We shall discuss how an investor might select an appropriate utility function after we examine a few more properties of utility functions and study some examples of their use.

**Example 11.1 (The venture capitalist)** Sybil, a venture capitalist, is considering two possible investment alternatives for the coming year. Her first alternative is to buy Treasury bills, which will give her a wealth of \$6M for sure. The second alternative has three possible outcomes. They will produce wealth levels \$10M, \$5M, and \$1M with corresponding probabilities of .2, .4, and .4. She decides to use the power utility  $U(x) = x^{1/2}$  to evaluate these alternatives (where  $x$  is in millions of dollars).

The first alternative has an expected utility of  $\sqrt{6} = 2.45$ . The second has an expected utility of  $.2 \times \sqrt{10} + .4 \times \sqrt{5} + .4 \times \sqrt{1} = .2 \times 3.16 + .4 \times 2.24 + .4 = 1.93$ . Hence the first alternative is preferred to the second.

There is good justification for using the expected value of a utility function as a basis for decision making. Indeed, the approach can be derived from a set of reasonable axioms that describe rational behavior (see references at end of chapter). Overall, this method has the merit of simplicity, good flexibility due to the possibility of selecting a variety of utility functions, and strong theoretical justification.

## Equivalent Utility Functions

Since a utility function is used to provide a ranking among alternatives, its actual numerical value (its cardinal value) has no real meaning. All that matters is how it

ranks alternatives when an expected utility is computed. It seems clear that a utility function can be modified in certain elementary ways without changing the rankings that it provides. We investigate this property here.

First, it is clear that the addition of a constant to a utility function does not affect its rankings. That is, if we use a utility function  $U(x)$  and then define the alternative function  $V(x) = U(x) + b$ , this new function provides exactly the same rankings as the original. This follows from the linearity of the expected value operation. Specifically,  $E[V(x)] = E[U(x) + b] = E[U(x)] + b$ . Hence the new expected utility values are equal to the old values plus the constant  $b$ . This addition does not change the rankings of various alternatives.

In a similar fashion it can be seen that the use of the function  $V(x) = aU(x)$  for a constant  $a > 0$  does not change the ranking because  $E[V(x)] = E[aU(x)] = aE[U(x)]$ .

In general, given a utility function  $U(x)$ , any function of the form

$$V(x) = aU(x) + b \quad (11.1)$$

with  $a > 0$  is a utility function **equivalent** to  $U(x)$ . Equivalent utility functions give identical rankings. [It can be shown that the transformation (11.1) is the only transformation that leaves the rankings of all random outcomes the same.] As an example, the utility function  $V(x) = \ln(cx^a)$  with  $a > 0$  is equivalent to the logarithmic utility function  $U(x) = \ln x$  because  $\ln(cx^a) = a \ln x + \ln c$ .

In practice, we recognize that a utility function can be changed to an equivalent one, and we may use this fact to scale a utility function conveniently.

## 11.3 Risk Aversion

The main purpose of a utility function is to provide a systematic way to rank alternatives that captures the principle of risk aversion. This is accomplished whenever the utility function is concave. We spell out this definition formally:

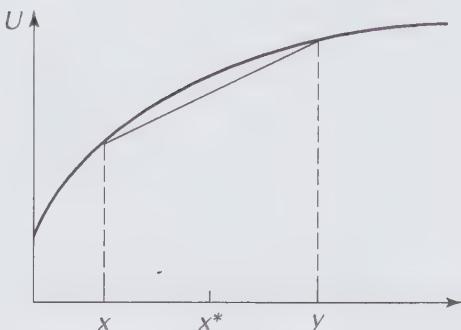
**Concave utility and risk aversion** A function  $U$  defined on an interval  $[a, b]$  of real numbers is said to be **concave** if for any  $\alpha$  with  $0 \leq \alpha \leq 1$  and any  $x$  and  $y$  in  $[a, b]$  there holds

$$U[\alpha x + (1 - \alpha)y] \geq \alpha U(x) + (1 - \alpha)U(y). \quad (11.2)$$

An (increasing) utility function  $U$  is said to be **risk averse** on  $[a, b]$  if it is concave on  $[a, b]$ . If  $U$  is concave everywhere, it is said to be **risk averse**.

This definition is illustrated in Figure 11.2. The figure shows a utility function that is concave. To check the concavity we take two arbitrary points  $x$  and  $y$  as shown, and any  $\alpha$ ,  $0 \leq \alpha \leq 1$ . The point  $x^* = \alpha x + (1 - \alpha)y$  is a weighted average of  $x$  and  $y$ , and hence  $x^*$  is between  $x$  and  $y$ . The value of the function at this point is greater than the value at  $x^*$  of the straight line connecting the function values at  $U(x)$  and  $U(y)$ . In general, the condition for concavity is that the straight line drawn between two points on the function must lie below (or on) the function itself. In simple terms, an increasing concave function has a slope that flattens for increasing values.

**FIGURE 11.2 Concavity and risk aversion.** The straight line connecting  $x$  and  $y$  lies below the function at any intermediate point. As a special case, a sure value of  $x^* = \frac{1}{2}x + \frac{1}{2}y$  is preferred to a 50–50 chance of  $x$  or  $y$ .



The same figure can be used to show how concavity of the utility function is related to risk aversion. Suppose that we have two alternatives for future wealth. The first is that we obtain either  $x$  or  $y$ , each with a probability of  $\frac{1}{2}$ . The second is that we obtain  $\frac{1}{2}x + \frac{1}{2}y$  with certainty. Suppose our utility function is the one shown in Figure 11.2. The expected utility of the first alternative (the 50–50 chance) is equal to the value of the straight line at the point  $x^* = \frac{1}{2}x + \frac{1}{2}y$ , because this is the weighting of the two utility values. The expected utility of the second option (the riskless one) is equal to the value of the function at the point  $x^* = \frac{1}{2}x + \frac{1}{2}y$ . This value is greater than that of the first alternative when the utility function is concave. Hence the sure wealth of  $\frac{1}{2}x + \frac{1}{2}y$  is preferred to a 50–50 chance of  $x$  or  $y$ . Both alternatives have the same expected value, but the one without risk is preferred.

A special case is the risk-neutral utility function  $U(x) = x$  [and its equivalent forms  $V(x) = ax + b$  with  $a > 0$ ]. This function is concave according to the preceding definition, but it is a limiting case. Strictly speaking, this function represents risk aversion of zero. Frequently we reserve the phrase *risk averse* for the case where  $U$  is *strictly concave*, which means that there is strict inequality in (11.2) whenever  $x \neq y$ .

**Example 11.2 (A coin toss)** As a specific example suppose that you face two options. The first is based on a toss of a coin—heads, you win \$10; tails, you win nothing. The second option is that you can have an amount  $M$  for certain. Your utility function for money is  $x - .04x^2$ . Let us evaluate these two alternatives. The first has expected utility  $E[U(x)] = \frac{1}{2}(10 - .04 \times 10^2) + \frac{1}{2}0 = 3$ . The second alternative has expected utility  $M - .04M^2$ . If  $M = 5$ , for example, then this value is 4, which is greater than the value of the first alternative. This means that you would favor the second alternative; that is, you would prefer to have \$5 for sure rather than a 50–50 chance of getting \$10 or nothing.

We can go a step further and determine what value of  $M$  would give the same utility as the first option. We solve  $M - .04M^2 = 3$ . This gives  $M = \$3.49$ . Hence you would be indifferent between getting \$3.49 for sure and having a 50–50 chance of getting \$10 or 0.

## Derivatives

We can relate important properties of a utility function to its derivatives. First,  $U(x)$  is strictly increasing with respect to  $x$  if  $U'(x) > 0$ . Second,  $U(x)$  is strictly concave with respect to  $x$  if  $U''(x) < 0$ . For example, consider the exponential utility function  $U(x) = -e^{-ax}$ . We find  $U'(x) = ae^{-ax} > 0$ , so  $U$  is increasing. Also,  $U''(x) = -a^2 e^{-ax} < 0$ , so  $U$  is concave.

## Risk Aversion Coefficients

The degree of risk aversion exhibited by a utility function is related to the magnitude of the bend in the function—the stronger the bend, the greater the risk aversion. This notion can be quantified in terms of the second derivative of the utility function.

The degree of risk aversion is formally defined by the **Arrow-Pratt absolute risk aversion coefficient**, which is

$$a(x) = -\frac{U''(x)}{U'(x)}.$$

The term  $U'(x)$  appears in the denominator to normalize the coefficient. With this normalization  $a(x)$  is the same for all equivalent utility functions. Basically, the coefficient function  $a(x)$  shows how risk aversion changes with the wealth level. For many individuals, risk aversion decreases as their wealth increases, reflecting the fact that they are willing to take more risk when they are financially secure.

As a specific example consider again the exponential utility function  $U(x) = -e^{-ax}$ . We have  $U'(x) = ae^{-ax}$  and  $U''(x) = -a^2 e^{-ax}$ . Therefore  $a(x) = a$ . In this case the risk aversion coefficient is constant for all  $x$ . If we make the same calculation for the equivalent utility function  $U(x) = 1 - be^{-ax}$ , we find that  $U'(x) = bae^{-ax}$  and  $U''(x) = -ba^2 e^{-ax}$ . So again  $a(x) = a$ .

As another example, consider the logarithmic utility function  $U(x) = \ln x$ . Here  $U'(x) = 1/x$  and  $U''(x) = -1/x^2$ . Therefore  $a(x) = 1/x$ ; and in this case, risk aversion decreases as wealth increases.

## Certainty Equivalent

Although the actual value of the expected utility of a random wealth variable is meaningless except in comparison with that of another alternative, there is a derived measure with units that do have intuitive meaning. This measure is the **certainty equivalent**.<sup>1</sup>

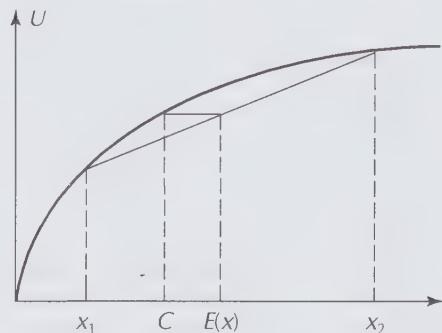
The certainty equivalent of a random wealth variable  $x$  is defined to be the amount of a certain (that is, risk-free) wealth that has a utility level equal to the

---

<sup>1</sup> This general concept of certainty equivalent is indirectly related to the concept with the same name used in Section 7.7.

**FIGURE 11.3 Certainty equivalent.**

The certainty equivalent is always less than the expected value for a risk-averse investor.



expected utility of  $x$ . In other words, the certainty equivalent  $C$  of a random wealth variable  $x$  is that value  $C$  satisfying

$$U(C) = E[U(x)].$$

The certainty equivalent of a random variable is the same for all equivalent utility functions and is measured in units of wealth.

As an example, consider the coin toss example discussed earlier. Our computation at the end of the example found that the certainty equivalent of the 50–50 chance of winning \$10 or \$0 is \$3.49 because that is the value that, if obtained with certainty, would have the same utility as the reward based on the outcome of the coin toss.

For a concave utility function it is always true that the certainty equivalent of a random outcome  $x$  is less than or equal to the expected value; that is,  $C \leq E(x)$ . Indeed, this inequality is another (equivalent) way to define risk aversion.

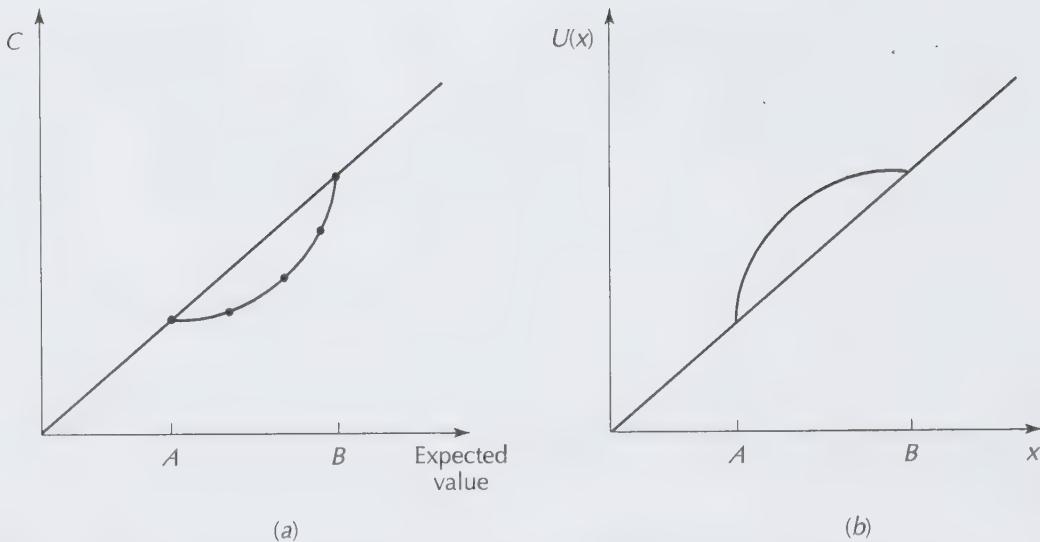
The certainty equivalent is illustrated in Figure 11.3 for the case of two outcomes  $x_1$  and  $x_2$ . The certainty equivalent is found by moving horizontally leftward from the point where the line between  $U(x_1)$  and  $U(x_2)$  intersects the vertical line drawn at  $E(x)$ .

## 11.4 Specification of Utility Functions\*

There are systematic procedures for assigning an appropriate utility function to an investor, some of which are quite elaborate. We outline a few general approaches in simple form.

### Direct Measurement of Utility

One way to measure an individual's utility function is to ask the individual to assign certainty equivalents to various risky alternatives. One particularly elegant way to organize this process is to select two fixed wealth values  $A$  and  $B$  as reference points. A lottery is then proposed that has outcome  $A$  with probability  $p$  and outcome  $B$  with probability  $1 - p$ . For various values of  $p$  the investor is asked how much certain



**FIGURE 11.4 Experimental determination of utility function.** (a) For lotteries that pay either  $A$  or  $B$  and have expected value  $e$ , a person is asked to state the certainty equivalent  $C$ . (b) Inverting this relation gives the utility function.

wealth  $C$  he or she would accept in place of the lottery.  $C$  will vary as  $p$  changes. Note that the values  $A$ ,  $B$ , and  $C$  are values for total wealth, not just increments based on a bet. A lottery with probability  $p$  has an expected value of  $e = pA + (1 - p)B$ . However, a risk-averse investor would accept less than this amount to avoid the risk of the lottery. Hence  $C < e$ .

The values of  $C$  reported by the investor for various  $p$ 's are plotted in Figure 11.4(a). The value of  $C$  is placed above the corresponding  $e$ . A curve is drawn through these points, giving a function  $C(e)$ . To define a utility function from this diagram, we normalize by setting  $U(A) = A$  and  $U(B) = B$  (which is legitimate because a utility function has two degrees of scaling freedom). With this normalization, the expected utility of the lottery is  $pU(A) + (1 - p)U(B) = pA + (1 - p)B$ , which is exactly the same as the expected value  $e$ . Therefore since  $C$  is defined so that  $U(C)$  is the expected utility of the lottery, we have the relation  $U(C) = e$ . Hence  $C = U^{-1}(e)$ , and thus the curve defined by  $C(e)$  is the inverse of the utility function. The utility function is obtained by flipping the axes to obtain the inverse function, as shown in Figure 11.4(b).<sup>2</sup>

**Example 11.3 (The venture capitalist)** Sybil, who has become a moderately successful venture capitalist, is anxious to make her utility function explicit. A consultant asks her to consider lotteries with outcomes of either \$1M or \$9M. She is asked to

<sup>2</sup> If different values of  $A$  and  $B$  are used, a new utility function is obtained, which is equivalent to the original one; that is, it is just a linear transformation of the original one. (See Exercise 5.)

**TABLE 11.1**  
**EXPECTED UTILITY VALUES AND CERTAINTY EQUIVALENTS**

<i>p</i>	0	.1	.2	.3	.4	.5	.6	.7	.8	.9	1
<i>e</i>	9	8.2	7.4	6.6	5.8	5	4.2	3.4	2.6	1.8	1
<i>C</i>	9	7.84	6.76	5.76	4.84	4	3.24	2.56	1.96	1.44	1

follow the direct procedure as the probability *p* of receiving \$1M varies. For a 50–50 chance of the two outcomes, the expected value is \$5M, but she assigns a certainty equivalent of \$4M. Other values she assigns are shown in Table 11.1.

The utility function is also shown in Table 11.1, since  $U(C) = e$ . (We just read from the bottom row up to the next row to evaluate *U*.) For example,  $U(4) = 5$ . However, the values of *C* in the table are not all whole numbers, so the table is not in the form that one would most desire. A new table of utility values could be constructed by interpolating in Table 11.1. For example (although perhaps not obviously),

$$U(2) = \frac{3.4(2.00 - 1.96) + 2.6(2.56 - 2.00)}{2.56 - 1.96} = 2.65.$$

## Parameter Families

Another simple method of assigning a utility function is to select a parameterized family of functions and then determine a suitable set of parameter values.

This technique is often carried out by assuming that the utility function is of the exponential form  $U(x) = -e^{-ax}$ . It is then only necessary to determine the parameter *a*, which is the risk aversion coefficient for this utility function. This parameter can be determined by evaluating a single lottery in certainty equivalent terms. For example, we might ask an investor how much he or she would accept in place of a lottery that offers a 50–50 chance of winning \$1 million or \$100,000. Suppose the investor felt that this was equivalent to a certain wealth of \$400,000. We then set

$$-e^{-400,000a} = -.5e^{-1,000,000a} - .5e^{-100,000a}.$$

We can solve this (by an iterative procedure) to obtain  $a = 1/\$623,426$ .

Many people prefer to use a logarithmic or power utility function, since these functions have the property that risk aversion decreases with wealth. Indeed, for the logarithmic utility, the risk aversion coefficient is  $a(x) = 1/x$ , and for the power utility function  $U(x) = \lambda x^\gamma$  the coefficient is  $a(x) = (1 - \gamma)/x$ . There are also good arguments based on the theory of Chapter 18, which suggest that these are appropriate utility functions for investors concerned with the long-term growth of their wealth.

A compromise, or composite, approach that is commonly used is to recognize that while utility is a function of total wealth, most investment decisions involve relatively small increments to that wealth. Hence if  $x_0$  is the initial wealth and *w* is the increment, the proper function is  $U(x_0 + w)$ . This is approximated by evaluating increments directly with an exponential utility function  $-e^{-aw}$ . However,

if we assume that the true utility function is  $\ln x$ , then we use  $a = 1/x_0$  in the exponential approximation.

**Example 11.4 (Curve fitting)** The tabular results of Example 11.3 (for the venture capitalist Sybil) can be expressed compactly by fitting a curve to the results. If we assume a power utility function, it will have the form  $U(x) = ax^\gamma + c$ . Our normalization requires

$$a + c = 1$$

$$a9^\gamma + c = 9.$$

Thus  $a = 8/(9^\gamma - 1)$  and  $c = (9^\gamma - 9)/(9^\gamma - 1)$ . Therefore it only remains to determine  $\lambda$ . We can find the best value to fit the values matching  $U(C)$  to  $e$  in Table 11.1. We find (using a spreadsheet optimizer) that, in fact,  $\gamma = \frac{1}{2}$  provides an excellent fit. Hence we set  $U(x) = 4\sqrt{x} - 3$ ; or as an equivalent form,  $V(x) = \sqrt{x}$ .

## Questionnaire Method

The risk aversion characteristics of an individual depends on the individual's feelings about risk, his or her current financial situation (such as net worth), the prospects for financial gains or requirements (such as college expenses), and the individual's age. One way, therefore, to attempt to deduce the appropriate risk factor and utility function for wealth increments is to administer a questionnaire such as the one shown in Figure 11.5. This gives a qualitative evaluation, and the results can be used to assign a specific function if desired.

In the questionnaire, note that section A concerns personal environment and objectives while section B is related to investment experience, and section C relates to the investors attitude towards risk. The questionnaire may influence suggested portfolios to a client. For example, complex financial products would not be suggested to inexperienced investors. Moreover, if an investor is inexperienced, we might want to treat answers to questions in section C with care.

## 11.5 Utility Functions and the Mean–Variance Criterion\*

The mean–variance criterion used in the Markowitz portfolio problem can be reconciled with the expected utility approach in either of two ways: (1) using a quadratic utility function, or (2) making the assumption that the random variables that characterize returns are normal (Gaussian) random variables. These two special cases are examined here.

### Quadratic Utility

The quadratic utility function can be defined as  $U(x) = ax - \frac{1}{2}bx^2$ , where  $a > 0$  and  $b \geq 0$ . This function is shown in Figure 11.6.

## RISK TOLERANCE PROFILE

### A. Investment Goal

1. What is your Investment horizon?
  - a) less than 20 years
  - b) between 20 and 40 years
  - c) more than 40 years
2. Do you want to access your investment to maintain your lifestyle?
  - a) No
  - b) Yes, but only from income
  - c) Yes
3. Which situation describes you best?
  - a) Single or married, no children
  - b) Family supporting dependents
  - c) Retired or close to retirement age

### B. Investment Experience

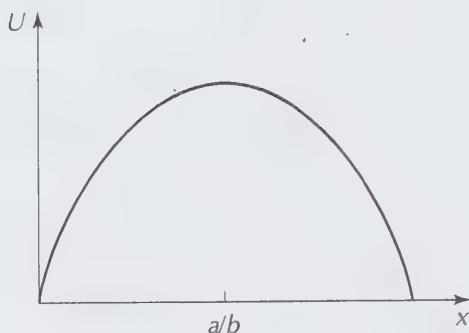
1. How long have you been making financial investments (with or without adviser)?
  - a) Less than 1 year or none
  - b) Between 1 and 5 years
  - c) More than 5 years
2. How would you describe your knowledge of financial markets?
  - a) None or little
  - b) Average
  - c) Expert

### C. Risk capacity and Preferences

1. What percentage of invested wealth would you need for future expenditures? \_\_\_\_ %
2. Losing which percentage of invested wealth would give you sleepless nights? \_\_\_\_ %

**FIGURE 11.5 Risk quiz.** An investor's attitude towards risk and toward type of investment might be inferred from responses to a questionnaire such as this one.

**FIGURE 11.6 Quadratic utility function.** This function is meaningful as a utility function only for  $x \leq a/b$ .



This utility function is really meaningful only in the range  $x \leq a/b$ , for it is in this range that the function is increasing. Note also that for  $b > 0$  the function is strictly concave everywhere and thus exhibits risk aversion.

We assume that all random variables of interest lie in the feasible range  $x \leq a/b$ ; that is, within the meaningful range of the quadratic utility function.

Suppose that a portfolio has a random wealth value of  $y$ . Using the expected utility criterion we evaluate the portfolio using the value

$$\begin{aligned} E[U(y)] &= E\left(ay - \frac{1}{2}by^2\right) \\ &= aE(y) - \frac{1}{2}bE(y^2) \\ &= aE(y) - \frac{1}{2}b[E(y)]^2 - \frac{1}{2}b\text{var}(y). \end{aligned}$$

The optimal portfolio is the one that maximizes this value with respect to all feasible choices of the random wealth variable  $y$ .

This can be seen to be equivalent to a mean-variance approach. First, for convenience, suppose that the initial wealth is 1. Then  $y$  corresponds exactly to the return  $R$ . Suppose also that the solution has an expected value  $E(y) = M$ . Then clearly,  $y$  must have minimum variance with respect to all feasible  $y$ 's with  $E(y) = M = 1 + m$  (where  $m$  is the mean rate of return). Since  $y = R$ , it follows that the solution must correspond to a mean-variance efficient point.

Different mean-variance efficient points are obtained by selecting different values for the parameters  $a$  and  $b$ . Likewise, if the initial wealth is not 1, a different factor is introduced.

## Normal Returns

When all returns are normal random variables, the mean-variance criterion is also equivalent to the expected utility approach for any risk-averse utility function. To deduce this, select a utility function  $U$ . Consider a random wealth variable  $y$  that is a normal random variable with mean value  $M$  and standard deviation  $\sigma$ . Since the probability distribution is completely defined by  $M$  and  $\sigma$ , it follows that the expected

utility is a function of  $M$  and  $\sigma$ ; that is,

$$E[U(y)] = f(M, \sigma).$$

(It may be impossible to determine the function  $f$  in closed form, but that does not matter.) If  $U$  is risk averse, then  $f(M, \sigma)$  will be increasing with respect to  $M$  and decreasing with respect to  $\sigma$ . Now suppose that the returns of all assets are normal random variables. Then (and this is the key property) any linear combination of these assets is a normal random variable, with some mean and standard deviation. (See Appendix A.) Hence any portfolio of these assets will have a return that is a normal random variable. The portfolio problem is therefore equivalent to the selection of that combination of assets that maximizes the function  $f(M, \sigma)$  with respect to all feasible combinations. For a risk-averse utility this again implies that the variance should be minimized for any given value of the mean. In other words, the solution must be mean-variance efficient. Therefore the mean-variance criterion is appropriate when all returns are normal random variables.

## 11.6 Linear Pricing

We now turn attention to a fundamental property of security pricing—namely, that of linearity. We shall find that this property has profound implications and by itself explains much of the theory developed in previous chapters. (The remaining sections of this chapter might best be read after completing Part 3.)

We formalize the definition of a **security** as a random payoff variable, say,  $d$ . The payoff is revealed and obtained at the end of the period. (The payoff can be thought of as a dividend, which justifies the use of the letter  $d$ .) Associated with a security is a price  $P$ . As an example, we can imagine a security that pays  $d = \$10$  if it rains tomorrow or  $d = -\$10$  if it is sunny, with zero initial price. (This would correspond to a \$10 bet that it will rain.) Or we could consider a share of XYZ stock whose value at the end of a year is unknown. The payoff  $d$  is that random value. The price is the current price of a share of XYZ.

### Type A Arbitrage

Linear pricing of securities follows from an assumption that the most basic form of arbitrage is not possible. We define this basic form of arbitrage as follows. If an investment produces an immediate positive reward with no future payoff (either positive or negative), that investment is said to be a **type A arbitrage**.

In other words, if you invest in a type A arbitrage, you obtain money immediately and never have to pay anything. You invest in a security that pays zero with certainty but has a negative price. It seems quite reasonable to assume that such things do not exist.

To see that linear pricing follows from the assumption that there is no possibility of type A arbitrage, suppose that  $d$  is a security with price  $P$ . Consider the security  $2d$  that always pays exactly twice what  $d$  pays. Suppose that its price were  $P' < 2P$ .

Then we could buy this double security at the reduced price, and then break it apart and sell the two halves at price  $P$  for each half. We would obtain a net profit of  $2P - P'$  and then have no further obligation, since we sold what we bought. We have an immediate profit, and hence have found a type A arbitrage. This argument can be reversed to show that the price of the double security cannot be greater than  $2P$ . The argument also can be extended to show that for any real number  $\alpha$  the price of  $\alpha d$  must be  $\alpha P$ .

Likewise, if  $d_1$  and  $d_2$  are securities with prices  $P_1$  and  $P_2$ , the price of the security  $d_1 + d_2$  must be  $P_1 + P_2$ . For if the price of  $d_1 + d_2$  were  $P' < P_1 + P_2$ , we could purchase the combined security for  $P'$ , then break it into  $d_1$  and  $d_2$  and sell these for  $P_1$  and  $P_2$ , respectively. As a result we would obtain a profit of  $P_1 + P_2 - P' > 0$ . As before, this argument can be reversed if  $P' > P_1 + P_2$ . Hence the price of  $d_1 + d_2$  must be  $P_1 + P_2$ . In general, therefore, the price of  $\alpha d_1 + \beta d_2$  must be equal to  $\alpha P_1 + \beta P_2$ . This is **linear pricing**.<sup>3</sup>

In addition to the absence of type A arbitrage, the preceding argument assumes an ideal functioning of the market: it assumes that securities can be arbitrarily divided into two pieces, and it assumes that there are no transaction costs. In practice these requirements are not met perfectly, but when dealing with large numbers of shares of traded securities in highly liquid markets, they are closely met.

## Portfolios

Suppose now that there are  $n$  securities  $d_1, d_2, \dots, d_n$ . A **portfolio** of these securities is represented by an  $n$ -dimensional vector  $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ . The  $i$ th component  $\theta_i$  represents the amount of security  $i$  in the portfolio. The payoff of the portfolio is the random variable

$$d = \sum_{i=1}^n \theta_i d_i.$$

Under the assumption of no type A arbitrage, the price of the portfolio  $\theta$  is found by linearity. Thus the total price is

$$P = \sum_{i=1}^n \theta_i P_i$$

which is a more general expression of linear pricing.

Recall that the CAPM formula in pricing form is linear.

## Type B Arbitrage

Another form of arbitrage can be identified. If an investment has nonpositive cost but has a positive probability of yielding a positive payoff and no probability of yielding a negative payoff, that investment is said to be a **type B arbitrage**.

---

<sup>3</sup> Linear pricing also follows from the **law of one price**: if  $d_1 = d_2$  then  $P_1 = P_2$ .

In other words, a type B arbitrage is a situation where an individual pays nothing (or a negative amount) and has a chance of getting something. An example would be a free lottery ticket—you pay nothing for the ticket, but have a chance of winning a prize. Clearly, such tickets are rare in securities markets.

The two types of arbitrage are distinguished only for clarity of the concepts involved. In further developments we shall usually assume that neither type A nor type B is possible, and we shall just say that there is **no arbitrage possibility**. However, we have shown that ruling out type A is all that is needed to establish linear pricing. Ruling out type B as well allows us to develop stronger relations, as shown in the next section.

## 11.7 Portfolio Choice

We are now prepared to put many of the earlier sections of this chapter together and consider the portfolio problem of an investor who uses an expected utility criterion to rank alternatives.

If  $x$  is a random variable, we write  $x \geq 0$  to indicate that the variable is never less than zero. We write  $x > 0$  to indicate that the variable is never less than zero and it is strictly positive with some positive probability.

Suppose that an investor has a strictly increasing utility function  $U$  and an initial wealth  $W$ . There are  $n$  securities  $d_1, d_2, \dots, d_n$ . The investor wishes to form a portfolio to maximize the expected utility of final wealth, say,  $x$ . We let the portfolio be defined by  $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ , which gives the amounts of the various securities. The investor's problem is

$$\text{maximize} \quad E[U(x)] \quad (11.3a)$$

$$\text{subject to} \quad \sum_{i=1}^n \theta_i d_i = x \quad (11.3b)$$

$$x \geq 0 \quad (11.3c)$$

$$\sum_{i=1}^n \theta_i P_i \leq W. \quad (11.3d)$$

This problem states that the investor must select a portfolio with total cost no greater than the initial wealth  $W$  (the last constraint), that the final wealth  $x$  is defined by the portfolio choice (the first constraint), that this final wealth must be nonnegative in every possible outcome (the second constraint), and that the investor wishes to maximize the expected utility of this final wealth.

We now show how this problem is connected to the arbitrage concepts.

**Portfolio choice theorem** *Suppose that  $U(x)$  is continuous, strictly concave, and strictly increasing toward infinity as  $x \rightarrow \infty$ . Suppose also that there is a portfolio  $\theta^0$  such that  $\sum_{i=1}^n \theta_i^0 d_i > 0$ . Then the optimal portfolio problem (11.3a) has a solution if and only if there is no arbitrage possibility.*

**Proof:** We shall only prove the *only if* portion of the theorem. Suppose that there is a type A arbitrage produced by a portfolio  $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ . Using this portfolio, it is possible to obtain additional initial wealth without affecting the final payoff. Hence arbitrary amounts of the portfolio  $\theta^0$  can be purchased. This implies that  $E[U(x)]$  does not have a maximum, because given a feasible portfolio, that portfolio can be supplemented by arbitrary amounts of  $\theta^0$  to increase  $E[U(x)]$ . If there is a type B arbitrage, it is possible to obtain (at zero or negative cost) an asset that has payoff  $\bar{x} > 0$  (with nonzero probability of being positive). We can acquire arbitrarily large amounts of this asset to increase  $E[U(x)]$  arbitrarily. Hence if there is a solution, there can be no type A or type B arbitrage. ■

We can go further than the preceding result on the existence of a solution and actually characterize the solution. We assume that there are no arbitrage opportunities and hence there is an optimal portfolio, which we denote by  $\theta^*$ . We also assume that the corresponding payoff  $x^* = \sum_{i=1}^n \theta_i^* d_i$  satisfies  $x^* > 0$ . We can immediately deduce that the inequality  $\sum_{i=1}^n \theta_i P_i \leq W$  will be met with equality at the solution; otherwise some positive fraction of the portfolio  $\theta^0$  (or  $\theta^*$ ) could be added to improve the result.

To derive the equations satisfied by the solution, we substitute  $x = \sum_{i=1}^n \theta_i d_i$  in the objective and ignore the constraint  $x \geq 0$  since we have assumed that it is satisfied by strict inequality. The problem therefore becomes

$$\begin{aligned} &\text{maximize} \quad E\left[U\left(\sum_{i=1}^n \theta_i d_i\right)\right] \\ &\text{subject to} \quad \sum_{i=1}^n \theta_i P_i = W. \end{aligned}$$

Assume  $U$  is continuously differentiable. Then by introducing a Lagrange multiplier  $\lambda$  for the constraint, and using  $x^* = \sum_{i=1}^n \theta_i^* d_i$  for the payoff of the optimal portfolio, the necessary conditions are found by differentiating the Lagrangian (see Appendix B)

$$L = E\left[U\left(\sum_{i=1}^n \theta_i d_i\right)\right] - \lambda \left(\sum_{i=1}^n \theta_i P_i - W\right)$$

with respect to each  $\theta_i$ . This gives

$$E[U'(x^*)d_i] = \lambda P_i \tag{11.4}$$

for  $i = 1, 2, \dots, n$ . This represents  $n$  equations. The original budget constraint  $\sum_{i=1}^n \theta_i P_i = W$  is one more equation. Altogether, therefore, there are  $n + 1$  equations for the  $n + 1$  unknowns  $\theta_1, \theta_2, \dots, \theta_n$  and  $\lambda$ . It can be shown that  $\lambda > 0$ .

These equations are very important because they serve two roles. First, and most obviously, they give enough equations to actually solve the optimal portfolio problem. An example of such a solution is given soon in Example 11.5. Second, since these equations are valid if there are no arbitrage opportunities, they provide a valuable characterization of prices under the assumption of no arbitrage. This use of the equations is explained in the next section.

If there is a risk-free asset with total return  $R$ , then equation (11.4) must apply when  $d_i = R$  and  $P_i = 1$ . Thus,

$$\lambda = E[U'(x^*)]R.$$

Substituting this value of  $\lambda$  in (11.4) yields

$$\frac{E[U'(x^*)d_i]}{R E[U'(x^*)]} = P_i.$$

Because of the importance of these equations, we now highlight them:

**Portfolio pricing equation** If  $x^* = \sum_{i=1}^n \theta_i^* d_i$ ,  $x^* > 0$ , is a solution to the optimal portfolio problem (11.3a), then

$$E[U'(x^*)d_i] = \lambda P_i \quad (11.5)$$

for  $i = 1, 2, \dots, n$ , where  $\lambda > 0$ . If there is a risk-free asset with return  $R$ , then

$$\frac{E[U'(x^*)d_i]}{R E[U'(x^*)]} = P_i \quad (11.6)$$

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

**Example 11.5 (A film venture)** An investor is considering the possibility of investing in a venture to produce an entertainment film. He has learned that such ventures are quite risky. In this particular case he has learned that there are essentially three possible outcomes, as shown in Table 11.2: (1) with probability .3 his investment will be multiplied by a factor of 3, (2) with probability .4 the factor will be 1, and (3) with probability .3 he will lose the entire investment. One of these outcomes will occur in 2 years. He also has the opportunity to earn a total of 20% risk free over this period. He wants to know whether he should invest money in the film venture; and if so, how much?

This is a simplification of a fairly realistic situation. The expected return is  $.3 \times 3 + .4 \times 1 + .3 \times 0 = 1.3$ , which is somewhat better than what can be obtained risk free. How much would you invest in such a venture? Think about it for a moment.

The investor decides to use  $U(x) = \ln x$  as a utility function. This is an excellent general choice (as will be explained in Chapter 18). His problem is to select amounts

**TABLE 11.2  
THE FILM VENTURE**

	Return	Probability
High success	3.0	0.3
Moderate success	1.0	0.4
Failure	0.0	0.3
Risk free	1.2	1.0

There are three possible outcomes with associated total returns and probabilities shown. There is also a risk-free opportunity with total return 1.2.

$\theta_1$  and  $\theta_2$  of the two available securities, the film venture and the risk-free opportunity, each of which has a unit price of 1. Hence his problem is to select  $(\theta_1, \theta_2)$  to solve

$$\text{maximize } [.3 \ln(3\theta_1 + 1.2\theta_2) + .4 \ln(\theta_1 + 1.2\theta_2) + .3 \ln(1.2\theta_2)]$$

$$\text{subject to } \theta_1 + \theta_2 = W.$$

The necessary conditions from equation (11.5), or by direct calculation, are

$$\frac{.9}{3\theta_1 + 1.2\theta_2} + \frac{.4}{\theta_1 + 1.2\theta_2} = \lambda$$

$$\frac{.36}{3\theta_1 + 1.2\theta_2} + \frac{.48}{\theta_1 + 1.2\theta_2} + \frac{.36}{1.2\theta_2} = \lambda.$$

These two equations, together with the constraint  $\theta_1 + \theta_2 = W$ , can be solved for the unknowns  $\theta_1$ ,  $\theta_2$ , and  $\lambda$ . (A quadratic equation must be solved.) The result is  $\theta_1 = .089W$ ,  $\theta_2 = .911W$ , and  $\lambda = 1/W$ . In other words, the investor should commit 8.9% of his wealth to this venture; the rest should be placed in the risk-free security.

**Example 11.6 (Residual rights)** While pondering the possibility of investing in the film venture of the previous example, our investor discovers that it is also possible to invest in film residuals, which have a large payoff if the film is highly successful. Each dollar invested in residual rights produces \$6 if the venture has high success and zero in the other two cases. Now what should the investor do?

He must solve the portfolio optimization problem again with this new information. There are now three securities: the original film venture, the risk-free alternative, and residual rights. He will purchase these in amounts  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$ , respectively. The necessary equations are

$$\frac{.9}{3\theta_1 + 1.2\theta_2 + 6\theta_3} + \frac{.4}{\theta_1 + 1.2\theta_2} = \lambda$$

$$\frac{.36}{3\theta_1 + 1.2\theta_2 + 6\theta_3} + \frac{.48}{\theta_1 + 1.2\theta_2} + \frac{.36}{1.2\theta_2} = \lambda$$

$$\frac{1.8}{3\theta_1 + 1.2\theta_2 + 6\theta_3} = \lambda.$$

In addition there is the wealth constraint  $\theta_1 + \theta_2 + \theta_3 = W$ . These equations have solution  $\theta_1 = -1.0W$ ,  $\theta_2 = 1.5W$ ,  $\theta_3 = .5W$ , and  $\lambda = 1/W$ . In other words, the investor should short the ordinary film venture by an amount equal to his total wealth in order to invest in the other two alternatives.

## 11.8 Arbitrage Bounds

Consider a market of securities in which there is no arbitrage opportunity. Now suppose a new asset with (random) payoff  $d$  is to be adjoined to the market at a certain price  $P_d$ . It is possible that this addition will produce an arbitrage opportunity in the now-expanded market. In general, this possibility will depend on the price

$P_d$ . We show that there is both an upper bound and a lower bound such that for  $P_d$  between these bounds there is no arbitrage.

**Proposition 11.1** *Assume there is no arbitrage possibility in the market. Let  $d$  be a new asset to be included in the market at a price  $P_d$ . Then there is an upper bound  $P_d^u$  and a lower bound  $P_d^l$  such that no arbitrage is possible for  $P_d \in (P_d^l, P_d^u)$ . Arbitrage is possible for  $P_d \notin [P_d^l, P_d^u]$ .*

**Proof:** Arbitrage is possible only if the new asset is included at a nonzero level. Suppose the level is negative. Without loss of generality we may use  $-1$ . An arbitrage consists of a marketed asset  $m$  and  $-d$  and satisfies

$$m - d \geq 0 \quad (11.7)$$

$$P_m - P_d \leq 0, \quad (11.8)$$

with at least one of the inequalities being strict. (Note that the first inequality must apply to all values that  $d$  might attain.) If a value of  $P_d$  satisfies these conditions, then clearly any larger value of  $P_d$  will also satisfy them. To find a lowest feasible  $P_d$  we select  $m$  to minimize  $P_m$  subject to satisfying the first condition of an arbitrage, so  $P_d$  can drop to this level. Indeed, the minimum value of  $P_d$  will equal the minimum feasible value of  $P_m$ . Hence<sup>4</sup>

$$P_d^l = \min_m \{P_m : m - d \leq 0\}.$$

A similar argument using  $+1$  as the level of  $d$  produces the lower bound  $P_d^u = \max_m \{P_m : m + d \geq 0\}$ . ■

**Example 11.7 (Easy coin flip)** Suppose someone offers a bet  $d$  on the result of a coin flip: If heads, the bet pays \$3; if tails, it pays zero. What are the price bounds for this bet? We assume that the market consists of a risk-free asset with 0% interest. In this case,  $d$  is either \$3 or \$0. For the upper bound we need  $\min_m \{P_m : m - d \geq 0\}$ , and we must use  $d = 3$ . The solution to the maximization is  $m = 3$ . It follows that  $P_d^u = 3$ . Likewise  $P_d^l = 0$ . Thus, any price in the interval  $(0, 3)$  will not cause an arbitrage. But a price outside  $[0, 3]$  will cause an arbitrage.<sup>5</sup>

## 11.9 Zero-Level Pricing

Suppose an investor has utility function  $U$ , which is continuous, strictly increasing, and strictly concave. The investor has wealth  $W > 0$ . The market has no arbitrage possibility. Suppose that the investor determines an optimal portfolio of marketed assets according to his or her utility function and wealth. Under the technical assumptions

<sup>4</sup> The unknowns can be taken to be the (finite) set of weights of basic market variables. Hence the minimum is achieved, and this is a **linear programming problem**.

<sup>5</sup> For this example the endpoints of the interval do not cause an arbitrage.

of the *Portfolio Choice Theorem*, there will be an optimal portfolio  $x^*$ , and we assume  $x^* > 0$ .

Now suppose a new asset  $d$ , outside the existing market, is introduced at some price  $P_d$  within the no-arbitrage interval  $(P_d^l, P_d^u)$ . If this is a relatively low price in the range, the investor may choose to modify the optimal portfolio by including some weight of  $d$ . Alternatively, if the price is relatively high, the investor may wish to short the asset. Under mild technical assumptions, there is a price  $P_d$  in the no-arbitrage interval such that a new optimal portfolio will include  $d$  only at the zero level. That is, there is no incentive to include the asset directly or by shorting. This  $P_d$  is the **zero-level price** of  $d$  for that investor, also known as the **indifference price** or **marginal price**.

**Example 11.8 (Zero coin flip)** Consider again the bet on the coin flip of Example 11.7, where the payoff  $d$  is \$3 or \$0. How much would you pay for this opportunity? First, we know that  $P_d \in [0, 3]$ . However, most likely you would decide that  $P = \$1.50$  might be most appropriate. It should in fact turn out that this price is your zero-level price.

In general, the zero-level price is different for different people because their utility functions and wealth levels differ. However, in the coin flip example, the result is the same for everyone; that is, all will agree. When that is the case, the zero-level price is said to be **universal**. Notice that the zero-level price of an asset already in the market is its market price, for that is the price used in the determination of the optimal portfolio, and at that price there is no incentive to change its weight.

This zero-level price can be found by applying pricing equation (11.6) for the case of the expanded market that includes  $d$  when we know that the optimal portfolio is again  $x^*$ . Thus we have

$$P_d = \frac{E[U'(x^*)d]}{R E[U'(x^*)]}. \quad (11.9)$$

As a special case, suppose that  $d$  is statistically independent of the market, as is the case for the coin flip bet, where the coin flip itself is independent of all assets in the market. Then the expected value in the numerator separates into  $E[U'(x^*)] E[d]$ , and, after canceling terms, we have  $P_d = E[d]/R$ , as expected; and it is valid for every investor. Also, the projection method in the CAPM framework can be considered as producing zero-level prices. See Exercise 12. As another case, if asset payoffs are jointly normal, the zero-level prices will be universal. Another important case is where the market is partially complete. See Exercise 13.

In some cases a zero-level price may be universal within a restriction. For example, people with logarithmic utility will all have the same zero-level price, regardless of wealth, so that price is universal for log optimizers.

A zero-level price is frequently the price obtained by theoretical methods. It shows how the new asset is related to the market and the existing portfolio. Presumably, if the actual price is higher (lower) than the zero-level price, one would buy (sell short) some amount of the asset.

**Example 11.9 (Zero-level rights)** Example 11.6 shows that the residual rights is a very attractive opportunity, since an investor would like to invest heavily in it by

shorting the original venture. Instead, let us find the zero-level price of the residual rights. We assume that the investor has found the optimal portfolio of Example 11.5. The zero-level price of residual rights is then

$$\left[ \frac{.3}{3\theta_1 + 1.2\theta_2} \right] 6 = 1.323.$$

The original price was 1.0. At the higher price, the investor will not modify the original portfolio by including an investment in residual rights.

## 11.10 Log-Optimal Pricing\*

The pricing formula

$$E[U'(x^*)d_i] = \lambda P_i, \quad i = 1, 2, \dots, n \quad (11.10)$$

is a general result with many important ramifications. It can be transformed to produce a variety of convenient special pricing formulas. This section presents one especially elegant version.

We shall choose  $U(x) = \ln x$  and  $W = 1$  as a special case to investigate. The final wealth variable  $x^*$  is then the one that is associated with the marketed portfolio that maximizes the expected logarithm of final wealth. In this special case we denote this  $x^*$  by  $R^*$ , since  $R^*$  is the return that is optimal for the logarithmic utility. We refer to  $R^*$  as the **log-optimal return**.

Now let  $d_i$  be the payoff of a security in the market. Since  $d \ln x/dx = 1/x$ , pricing equation (11.10) becomes

$$E\left(\frac{d_i}{R^*}\right) = \lambda P_i \quad (11.11)$$

for all  $i$ . Since this is valid for every security  $i$ , it is, by linearity, valid for the log-optimal portfolio itself. This portfolio has price 1, and therefore we find that

$$1 = E\left(\frac{R^*}{R^*}\right) = \lambda.$$

Thus we have found the value of  $\lambda$  for this case.

If there is a risk-free asset, portfolio pricing equation (11.10) is valid for it as well. The risk-free asset has a payoff identically equal to 1 and price  $1/R$ , where  $R$  is the total risk-free return. Hence we find

$$E(1/R^*) = 1/R.$$

Therefore we know that the expected value of  $1/R^*$  is equal to  $1/R$ .

Using the value of  $\lambda = 1$ , pricing equation (11.11) becomes

$$P_i = E\left(\frac{d_i}{R^*}\right).$$

Since this is true for any security  $i$ , it is, by linearity, also true for any portfolio.

We now recognize that this same formula can be applied to any asset with a payoff  $d$ . This will produce that zero-level price of  $d$  for a log investor. Thus we have the following result.

**Log-optimal pricing** *The log-optimal zero-level price  $P$  of any asset with payoff  $d$  is*

$$P = E\left(\frac{d}{R^*}\right), \quad (11.12)$$

where  $R^*$  is the return on the log-optimal portfolio of market securities.

Isn't this a simple and easily remembered result? The formula looks very similar to the expression  $P = d/R$  that would hold in the case where  $d$  is deterministic. In the random case we just substitute  $R^*$  for  $R$  and put an expected value in front. If  $d$  happens to be deterministic, this more general result reduces to the simple one because  $E(1/R^*) = 1/R$ .

**Example 11.10 (Film variations)** Suppose that a new security is proposed with payoffs that depend only on the possible outcomes of the film venture. For example, one might propose an investment that paid back something even if the venture was a failure. A general security of this type will have payoffs  $d^1, d^2$ , and  $d^3$ , corresponding to high success, moderate success, and failure, respectively. We can find the appropriate price of such a security by using the log-optimal portfolio that we calculated in Example 11.6.

Note that we cannot use the simple log-optimal portfolio of the first film venture example, because that example only considered the film venture and the risk-free security. If a new security were a combination of those two, then we could use the simple log-optimal portfolio for pricing. But if the new security is a general one, we must use the log-optimal portfolio of the second example, since it includes a complete set of three securities for the three possibilities. Any new security will be a combination of these three.

The log-optimal portfolio has the following return:

	High Success	Moderate Success	Failure
$R^*$	1.8	.8	1.8

These returns are calculated from the  $\theta_i$ 's found in the residual rights example. For example, under high success  $R^* = -1.0 \times 3 + 1.5 \times 1.2 + .5 \times 6 = 1.8$ .

The value of a security with payoffs  $d^1, d^2, d^3$  is  $E(d/R^*)$ , which is

$$P = .3 \frac{d^1}{1.8} + .4 \frac{d^2}{8} + .3 \frac{d^3}{1.8}.$$

You can try this on the three securities we have used before; their prices should all turn out to be 1. For example, for the original venture,  $P = .3 \frac{3}{1.8} + .4 \frac{1}{8} = \frac{1}{2} + \frac{1}{2} = 1$ .

We shall return to this log-optimal pricing equation in Chapter 18. For the moment we may regard it simply as a special version of the general pricing equation—the version obtained by using  $\ln x$  as the utility function.

Remember what is happening here. The prices of the original securities were used to find  $x^*$ . Now we use  $x^*$  to find those prices again. However, since pricing is linear, we can find the price of any security that is a linear combination of the original ones by the same formula.

What about a new security  $d$  that is not a linear combination of the original ones? We could enter it into the pricing equation as well, but the price obtained this way may not be correct. The formula is valid only for the securities used to derive it, or for a linear combination of those original securities.

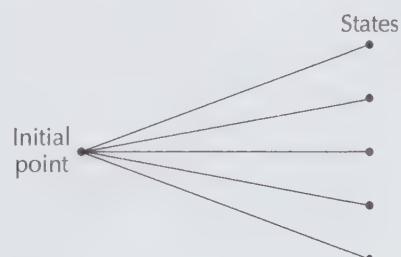
## 11.11 Finite State Models

Suppose that there are a finite number of possible **states** that describe the possible outcomes of a specific investment situation (see Figure 11.7). At the initial time it is known only that one of these states will occur. At the end of the period, one specific state will be revealed. Sometimes states describe certain physical phenomena. For example, we might define two weather states for tomorrow: sunny and rainy. We do not know today which of these will occur, but tomorrow this uncertainty will be resolved. Or, as another example, the states may correspond to economic events, as in the film venture example, which has the three possible states of high success, moderate success, and failure. Normally we index the possible states by numbers  $\{1, 2, \dots, S\}$ .

States define uncertainty in a very basic manner. It is not even necessary to introduce probabilities of the states, although this will be done later. Indeed, one of the main points of this section is that a great deal can be said without reference to probabilities. In an important sense, probabilities are irrelevant for pricing relations.

A **security** is defined within the context of states as a set of payoffs—one payoff for each possible state (again without reference to probabilities). Hence a security is represented by a vector of the form  $d = \langle d^1, d^2, \dots, d^S \rangle$ . We use the notation  $\langle \rangle$  to denote vectors whose components are state payoffs. In this case, the component  $d^s$ ,  $s = 1, 2, \dots, S$ , represents the payoff that is obtained if state  $s$  occurs. As before, associated with a security is a price  $P$ . Our earlier example, at the beginning of

**FIGURE 11.7 States.** States represent uncertainty in a simple but effective manner.



Section 11.6, of a security that pays \$10 if it rains tomorrow and -\$10 if it is sunny (with zero price), works here as well; and it is not necessary to specify probabilities. This security is represented as  $\langle 10, -10 \rangle$ .

## Completeness

Corresponding to the set of states, typically several securities are available. Security  $k$  may be written as  $d_k = \langle d_k^1, d_k^2, \dots, d_k^S \rangle$ . The set of available securities defines a market. As usual, a portfolio is a linear combination of securities, which forms a new security that is also in the market. If the market contains  $S$  linearly independent securities, the market is said to be **complete**. In this case, any vector of  $S$  components can be constructed to be the payoff of a security. In general, if there are fewer than  $S$  independent securities, the market is said to be **incomplete**.

## State Prices

A special form of security is one that has a payoff in only one state. Indeed, we can define the  $S$  **elementary state securities**  $e_s = \langle 0, 0, \dots, 0, 1, 0, \dots, 0 \rangle$ , where the 1 is the component  $s$  for  $s = 1, 2, \dots, S$ . If such a security exists, we denote its price by  $\psi_s$ .

When a complete set of state securities exists (one for each state), it is easy to determine the price of any other security. The security  $d = \langle d^1, d^2, \dots, d^S \rangle$  can be expressed as a combination of the elementary state securities as  $d = \sum_{s=1}^S d^s e_s$ , and hence by the linearity of pricing, the price of  $d$  must be

$$P = \sum_{s=1}^S d^s \psi_s. \quad (11.13)$$

If elementary state securities do not exist, it may be possible to construct them artificially by combining securities that do exist. For example, in a two-state world, if  $\langle 1, 1 \rangle$  and  $\langle 1, -1 \rangle$  exist, then one-half the sum of these two securities is equivalent to the first elementary state security  $\langle 1, 0 \rangle$ .

## Positive State Prices

If a complete set of elementary securities exists or can be constructed as a combination of existing securities, it is important that their prices be positive. Otherwise there would be an arbitrage opportunity. To see this, suppose an elementary state security  $e_s$  had a zero or negative price. That security would then present the possibility of obtaining something (a payoff of 1 if the state  $s$  occurs) for nonpositive cost. This is type B arbitrage. So if elementary state securities actually exist or can be constructed as combinations of other securities, their prices must be positive to avoid arbitrage.

Actually, the condition of no arbitrage possibility is equivalent to the existence of positive state prices as established by the following theorem:

**Positive state prices theorem** A set of positive state prices exists if and only if there are no arbitrage opportunities.

**Proof:** Suppose first that there are positive state prices. Then it is clear that no arbitrage is possible. To see this, suppose a security  $d$  can be constructed with  $d \geq 0$ . We have  $d = \langle d^1, d^2, \dots, d^S \rangle$  with  $d^s \geq 0$  for each  $s = 1, 2, \dots, S$ . The price of  $d$  is  $P = \sum_{s=1}^S \psi_s d^s$ , which since  $\psi_s > 0$  for all  $s$ , gives  $P \geq 0$ . Indeed  $P > 0$  if  $d \neq 0$  and  $P = 0$  if  $d = 0$ . Hence there is no arbitrage possibility.

To prove the converse, we assume that there are no arbitrage opportunities, and we make use of the result on the portfolio choice problem of Section 11.7. This proof requires some additional assumptions. (A more general proof is outlined in Exercise 14.) We assume there is a portfolio  $\theta^0$  such that  $\sum_{i=1}^n \theta_i^0 d_i > 0$ . We assign positive probabilities  $p_s, s = 1, 2, \dots, S$ , to the states arbitrarily, with  $\sum_{s=1}^S p_s = 1$ , and we select a strictly increasing utility function  $U$ . Since there is no arbitrage, there is, by the portfolio choice theorem of Section 11.7, a solution to the optimal portfolio choice problem. We assume that the optimal payoff has  $x^* > 0$ . The necessary conditions (11.4) show that for any security  $d$  with price  $P$ ,

$$E[U'(x^*)d] = \lambda P, \quad (11.14)$$

where  $x^*$  is the (random) payoff of the optimal portfolio and  $\lambda > 0$  is the Lagrange multiplier.

If we expand this equation to show the details of the expected value operation, we find

$$P = \frac{1}{\lambda} \sum_{s=1}^S p_s U'(x^*)^s d^s,$$

where  $U'(x^*)^s$  is the value of  $U'(x^*)$  in state  $s$ .

Now we define

$$\psi_s = \frac{p_s U'(x^*)^s}{\lambda}. \quad (11.15)$$

We see that  $\psi_s > 0$  because  $p_s > 0$ ,  $U'(x^*)^s > 0$ , and  $\lambda > 0$ . We also have

$$P = \sum_{s=1}^S \psi_s d^s$$

showing that the  $\psi_s$ 's are state prices. They are all positive. ■

Note that the theorem says that such positive prices exist—it does not say that they are unique. If there are more states than securities, there may be many different ways to assign state prices that are consistent with the prices of the existing securities. The theorem only says that for one of these ways the state prices are positive.

**Example 11.11 (The plain film venture)** Consider again the original film venture. There are three states but only two securities: the venture itself and the riskless security. Hence state prices are not unique.

We can find a set of positive state prices by using equation (11.15) and the values of the  $\theta_i$ 's and  $\lambda = 1$  found in Example 11.5 (with  $W = 1$ ). We have

$$\psi_1 = \frac{.3}{3\theta_1 + 1.2\theta_2} = .221$$

$$\psi_2 = \frac{.4}{\theta_1 + 1.2\theta_2} = .338$$

$$\psi_3 = \frac{.3}{1.2\theta_2} = .274.$$

These state prices can be used only to price combinations of the original two securities. They could not be applied, for example, to the purchase of residual rights. To check the price of the original venture we have  $P = 3 \times .221 + .338 = 1$ , as it should be.

**Example 11.12 (Expanded film venture)** Now consider the film venture with three available securities, as discussed in Example 11.6, which introduces residual rights. Since there are three states and three securities, the state prices are unique. Indeed we may find the state prices by setting the price of the three securities to 1, obtaining

$$\begin{aligned} 3\psi_1 + \psi_2 &= 1 \\ 1.2\psi_1 + 1.2\psi_2 + 1.2\psi_3 &= 1 \\ 6\psi_1 &= 1. \end{aligned}$$

This system has the solution

$$\psi_1 = \frac{1}{6}, \quad \psi_2 = \frac{1}{2}, \quad \psi_3 = \frac{1}{6}.$$

Therefore the price of a security with payoff  $\langle d^1, d^2, d^3 \rangle$  is

$$P = \frac{1}{6}d^1 + \frac{1}{2}d^2 + \frac{1}{6}d^3.$$

Note also that these state prices, although different from those of the preceding example, give the same values for prices of securities that are combinations of just the two in the original film venture. For example, the price of the basic venture itself is  $P = \frac{3}{6} + \frac{1}{2} = 1$ .

## 11.12 Risk-Neutral Pricing

Suppose there are positive state prices  $\psi_s, s = 1, 2, \dots, S$ . Then the price of any security  $d = \langle d^1, d^2, \dots, d^S \rangle$  can be found from

$$P = \sum_{s=1}^S d^s \psi_s.$$

We now normalize these state prices so that they sum to 1. Hence we let  $\psi_0 = \sum_{s=1}^S \psi_s$ , and let  $q_s = \psi_s / \psi_0$ . We can then write the pricing formula as

$$P = \psi_0 \sum_{s=1}^S q_s d^s. \quad (11.16)$$

The quantities  $q_s$ ,  $s = 1, 2, \dots, S$ , can be thought of as (artificial) probabilities, since they are positive and sum to 1. Using these as probabilities, we can write the pricing formula as

$$P = \psi_0 \hat{E}(d) \quad (11.17)$$

where  $\hat{E}$  denotes expectation with respect to the artificial probabilities  $q_s$ .

The value  $\psi_0$  has a useful interpretation. Since  $\psi_0 = \sum_{s=1}^S \psi_s$ , we see that  $\psi_0$  is the price of the security  $(1, 1, \dots, 1)$  that pays 1 in every state—a risk-free bond. By definition, its price is  $1/R$ , where  $R$  is the risk-free return. Thus we can write the pricing formula as

$$P = \frac{1}{R} \hat{E}(d). \quad (11.18)$$

This equation states that the price of a security is equal to the discounted expected value of its payoff, under the artificial probabilities. We term this **risk-neutral pricing** since it is exactly the formula that we would use if the  $q_s$ 's were real probabilities and we used a risk-neutral utility function (that is, the linear utility function). We also refer to the  $q_s$ 's as **risk-neutral probabilities**.

This artifice is deceptive in its simplicity; we shall find in the coming chapters that it has profound consequences. In fact a major portion of Part 3 is elaboration of this simple idea. Here are three ways to find the risk-neutral probabilities  $q_s$ :

- (a) The risk-neutral probabilities can be found from positive state prices by multiplying those prices by the risk-free rate. This is how we defined the risk-neutral probabilities at the beginning of this section.
- (b) If the positive state prices were found from a portfolio problem and there is a risk-free asset, we can use equation (11.6) to define

$$q_s = \frac{p_s U'(x^*)^s}{\sum_{t=1}^S p_t U'(x^*)^t}. \quad (11.19)$$

This formula will be useful in our later work.

- (c) If there are  $n$  states and at least  $n$  independent securities with known prices, and no arbitrage possibility, then the risk-neutral probabilities can be found directly by solving the system of equations

$$P_i = \frac{1}{R} \sum_{s=1}^S q_s d_i^s, \quad i = 1, 2, \dots, n$$

for the  $n$  unknown  $q_s$ 's.

**Example 11.13 (The film venture)** We found the state prices of the full film venture (with three securities) to be

$$\psi_1 = \frac{1}{6}, \quad \psi_2 = \frac{1}{2}, \quad \psi_3 = \frac{1}{6}.$$

Multiplying these by the risk-free rate 1.2, we obtain the risk-neutral probabilities

$$q_1 = .2, \quad q_2 = .6, \quad q_3 = .2.$$

Hence the price of a security with payoff  $\langle d^1, d^2, d^3 \rangle$  is

$$P = \frac{.2d^1 + .6d_2 + .2d_3}{1.2}.$$

Here again, this pricing formula is valid only for the original securities or linear combinations of those securities. The risk-neutral probabilities were derived explicitly to price the original securities.

The risk-neutral pricing result can be extended to the general situation that does not assume that there are a finite number of states. (See Exercise 17.)

## 11.13 Summary

This chapter is devoted to general theory, and hence it is somewhat more abstract than other chapters, but the tools presented are quite powerful. The chapter should be reviewed after reading Part 3 and again after reading Part 4.

The first part of the chapter presents the basics of expected utility theory. Utility functions account for risk aversion in financial decision making, and provide a more general and more useful approach than does the mean-variance framework. In this new approach, an uncertain final wealth level is evaluated by computing the expected value of the utility of the wealth. One random wealth level is preferred to another if the expected utility of the first is greater than that of the second. Often the utility function is expressed in analytic form. Commonly used functions are: exponential, logarithmic, power, and quadratic. A utility function  $U(x)$  can be transformed to  $V(x) = aU(x) + b$  with  $a > 0$ , and the new function  $V$  is equivalent to  $U$  for decision-making purposes.

It is generally assumed that a utility function is increasing, since more wealth is preferred to less. A utility function exhibits risk aversion if it is concave. If the utility function has derivatives and is both strictly increasing and strictly concave, then  $U'(x) > 0$  and  $U''(x) < 0$ .

Corresponding to a random wealth level, there is a number  $C$ , called the certainty equivalent of that random wealth. The certainty equivalent is the minimum (nonrandom) amount that an investor with utility function  $U$  would accept in place of the random wealth under consideration. The value  $C$  is defined such that  $U(C)$  is equal to the expected utility due to the random wealth level.

In order to use the utility function approach, an appropriate utility function must be selected. One way to make this selection is to assess the certain equivalents of various lotteries, and then work backward to find the underlying utility function that would assign those certain equivalent values.

Frequently the utility function is assumed to be either the exponential form  $-e^{-ax}$  with  $a$  approximately equal to the reciprocal of total wealth, the logarithmic form  $\ln x$ , or a power form  $\gamma x^\lambda$  with  $\lambda < 1$  but close to 0. The parameters of the function are either fit to lottery responses or deduced from the answers to a series of questions about an investor's financial situation and attitudes toward risk.

The second part of the chapter presents the outline of a general theory of linear pricing. In perfect markets (without transactions costs and with the possibility of buying or selling any amount of each security), security prices must be linear—meaning that the price of a bundle of securities must equal the sum of the prices of the component securities in the bundle—otherwise there is an arbitrage opportunity.

Two types of arbitrage are distinguished in the chapter: type A, which rules out the possibility of obtaining something for nothing—right now; and type B, which rules out the possibility of obtaining a chance for something later—at no cost now. Ruling out type A arbitrage leads to linear pricing. Ruling out both types A and B implies that the problem of finding the portfolio that maximizes the expected utility has a well-defined solution.

An investor with a given utility function and wealth can form an optimal portfolio of the marketed securities, provided no arbitrage is possible. This construction produces a linear pricing formula that will correctly price all marketed securities. If a new asset  $d$  becomes available, then there is a range of prices for which introduction of the asset at that particular price will not introduce an arbitrage in the expanded market. A given price in this band may, on reoptimization, cause the new asset to be included in the investor's portfolio at some level, long or short. There is a particular price  $P_d$  such that on reoptimization the optimal portfolio does not change. The new asset enters only at the zero level. This price is termed the **zero-level price** of the asset, and it can be found by extending the linear pricing rule for securities to include the new asset as well. The zero-level price is said to be **universal** if it is independent of the investor's utility function and wealth.

The optimal portfolio problem can be used to price realistic investments (such as the film venture). Furthermore, the necessary conditions of this general problem can be used in a backward fashion to express a security price as an expected value. Different choices of utility functions lead to different pricing formulas, but all of them are equivalent when applied to securities that are linear combinations of those considered in the original optimal portfolio problem. Utility functions that lead to especially convenient pricing equations include quadratic functions (which lead to the CAPM formula), and exponential, power, and logarithmic functions.

Insight and practical advantage can be derived from the use of finite state models. In these models it is useful to introduce the concept of state prices. A set of positive state prices consistent with the securities under consideration exists if and only if there are no arbitrage opportunities. One way to find a set of positive state prices is to solve the optimal portfolio problem. The state prices are determined directly by the resulting optimal portfolio.

A concept of major significance is that of risk-neutral pricing of securities. By introducing certain artificial probabilities, the pricing formula can be written as  $P = \hat{E}(d)/R$ , where  $R$  is the return of the riskless asset and  $\hat{E}$  denotes expectation with

respect to the artificial (risk-neutral) probabilities. A set of risk-neutral probabilities can be found by multiplying the state prices by the total return  $R$  of the risk-free asset.

## Exercises

1. (Certainty equivalent) An investor has utility function  $U(x) = x^{1/4}$  for salary. He has a new job offer which pays \$80,000 with a bonus. The bonus will be \$0, \$10,000, \$20,000, \$30,000, \$40,000, \$50,000, or \$60,000, each with equal probability. What is the certainty equivalent of this job offer?
2. (Wealth independence) Suppose an investor has exponential utility function  $U(x) = -e^{-ax}$  and an initial wealth level of  $W$ . The investor is faced with an opportunity to invest an amount  $w \leq W$  and obtain a random payoff  $x$ . Show that his evaluation of this incremental investment is independent of  $W$ .
3. (Risk aversion invariance) Suppose  $U(x)$  is a utility function with Arrow-Pratt risk aversion coefficient  $a(x)$ . Let  $V(x) = c + bU(x)$ . What is the risk aversion coefficient of  $V$ ?
4. (Relative risk aversion) The Arrow-Pratt relative risk aversion coefficient is

$$\mu(x) = \frac{xU''(x)}{U'(x)}.$$

Show that the utility functions  $U(x) = \ln x$  and  $U(x) = \gamma x^\gamma$  have constant relative risk aversion coefficients.

5. (Equivalency) A young woman uses the first procedure described in Section 11.4 to deduce her utility function  $U(x)$  over the range  $A \leq x \leq B$ . She uses the normalization  $U(A) = A$ ,  $U(B) = B$ . To check her result, she repeats the whole procedure over the range  $A' \leq x \leq B'$ , where  $A < A' < B' < B$ . The result is a utility function  $V(x)$ , with  $V(A') = A'$ ,  $V(B') = B'$ . If the results are consistent,  $U$  and  $V$  should be equivalent; that is,  $V(x) = aU(x) + b$  for some  $a > 0$  and  $b$ . Find  $a$  and  $b$ .

6. (HARA  $\diamond$ ) The HARA (for hyperbolic absolute risk aversion) class of utility functions is defined by

$$U(x) = \frac{1-\gamma}{\gamma} \left( \frac{ax}{1-\gamma} + b \right)^\gamma, \quad b > 0.$$

The functions are defined for those values of  $x$  where the term in parentheses is nonnegative. Show how the parameters  $\gamma$ ,  $a$ , and  $b$  can be chosen to obtain the following special cases (or an equivalent form).

- (a) Linear or risk neutral:  $U(x) = x$
- (b) Quadratic:  $U(x) = x - \frac{1}{2}cx^2$
- (c) Exponential:  $U(x) = -e^{-ax}$  [Try  $\gamma = -\infty$ .]
- (d) Power:  $U(x) = cx^\gamma$
- (e) Logarithmic:  $U(x) = \ln x$  [Try  $U(x) = (1-\gamma)^{1-\gamma}((x^\gamma - 1)/\gamma)$ .]

Show that the Arrow-Pratt risk aversion coefficient is of the form  $1/(cx + d)$ .

7. (The venture capitalist) A venture capitalist with a utility function  $U(x) = \sqrt{x}$  carried out the procedure of Example 11.3. Find an analytical expression for  $C$  as a function of  $e$ ,

and for  $e$  as a function of  $C$ . Do the values in Table 11.1 of the example agree with these expressions?

8. (Certainty approximation  $\diamond$ ) There is a useful approximation to the certainty equivalent that is easy to derive. A second-order expansion near  $\bar{x} = E(x)$  gives

$$U(x) \approx U(\bar{x}) + U'(\bar{x})(x - \bar{x}) + \frac{1}{2}U''(\bar{x})(x - \bar{x})^2.$$

Hence,

$$E[U(x)] \approx U(\bar{x}) + \frac{1}{2}U''(\bar{x})\text{var}(x).$$

On the other hand, if we let  $c$  denote the certainty equivalent and assume it is close to  $\bar{x}$ , we can use the first-order expansion

$$U(c) \approx U(\bar{x}) + U'(\bar{x})(c - \bar{x}).$$

Using these approximations, show that

$$c \approx \bar{x} + \frac{U''(\bar{x})}{U'(\bar{x})}\text{var}(x).$$

9. (Quadratic mean–variance) An investor with unit wealth maximizes the expected value of the utility function  $U(x) = ax - bx^2/2$  and obtains a mean–variance efficient portfolio. A friend of his with wealth  $W$  and the same utility function does the same calculation, but gets a different portfolio return. However, changing  $b$  to  $b'$  does yield the same result. What is the value of  $b'$ ?

10. (Portfolio optimization) Suppose an investor has utility function  $U$ . There are  $n$  risky assets with rates of return  $r_i$ ,  $i = 1, 2, \dots, n$ , and one risk-free asset with rate of return  $r_f$ . The investor has initial wealth  $W_0$ . Suppose that the optimal portfolio for this investor has (random) payoff  $x^*$ . Show that

$$E[U'(x^*)(r_i - r_f)] = 0$$

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

11. (Money-back guarantee) The promoter of the film venture offers a new investment designed to attract reluctant investors. One unit of this new investment has a payoff of three times the original investment if the venture is highly successful, and it refunds the original investment otherwise. Assuming that the other three investment alternatives described in Example 11.6 are also available, what is the price of this money-back guaranteed investment?

12. (Universal projection) Explain that in a CAPM framework, the projection price of a new asset is a universal zero-level price for all mean–variance investors.

13. (Universal partial) A market is **partially complete** if for any  $m$  in the market and any real function  $f$ , then  $f(m)$  is also in the market. In addition, there are other assets that are not traded.

As an example, suppose the market is made up of three states and two assets, A and B, as shown in the first four columns of the table. The market is not complete, since there are three states but only two assets. However, note that for each asset A or B, the value for state 2 is the same as for state 3. It follows that any function of these two assets will also have this property, and a combination of A and B can duplicate the function.

State	Prob.	A	B	$U'$	$d$
$s_1$	1/4	1	0	$u_1$	1
$s_2$	1/2	0	1	$u_2$	2
$s_3$	1/4	0	1	$u_2$	3

The fifth column shows the components of  $U'(x^*)$ , where  $x^*$  is the optimal portfolio. Again the entries for states 2 and 3 are equal. We assume things are scaled, so the price of an asset  $x$  in the market is  $P_x = E[U'(x^*)x]$ . The vector of  $U'(x^*)$  is unique since it can price all assets in the market.

The asset  $d = (1, 2, 3)$  is not in the marketed space. We evaluate its price as  $E[U'(x^*)d] = \frac{1}{4}u_1 + \frac{1}{2}2u_2 + \frac{1}{4}3u_2$ . Assuming the prices of A and of B are each \$1.00, find the zero-level price of  $d$ , and argue that it is universal.

14. (General positive state prices result  $\diamond$ ) The following is a general result from matrix theory: Let  $\mathbf{A}$  be an  $m \times n$  matrix. Suppose that the equation  $\mathbf{Ax} = \mathbf{p}$  can achieve no  $\mathbf{p} \geq \mathbf{0}$  except  $\mathbf{p} = \mathbf{0}$ . Then there is a vector  $\mathbf{y} > \mathbf{0}$  with  $\mathbf{A}^T \mathbf{y} = \mathbf{0}$ . Use this result to show that if there is no arbitrage, there are positive state prices; that is, prove the positive state price theorem in Section 11.9. [Hint: If there are  $S$  states and  $N$  securities, let  $\mathbf{A}$  be an appropriate  $(S+1) \times N$  matrix.]
15. (Quadratic pricing  $\diamond$ ) Suppose an investor uses the quadratic utility function  $U(x) = x - \frac{1}{2}cx^2$ . Suppose there are  $n$  risky assets and one risk-free asset with total return  $R$ . Let  $R_M$  be the total return on the optimal portfolio of risky assets. Show that the expected return of any asset  $i$  is given by the formula
- $$\bar{R}_i - R = \beta_i(\bar{R}_M - R)$$
- where  $\beta_i = \text{cov}(R_M, R_i)/\sigma_M^2$ . [Hint: Use Exercise 10. Apply the result to  $R_M$  itself.]
16. (At the track) At the horse races one Saturday afternoon Gavin Jones studies the racing form and concludes that the horse No Arbitrage has a 25% chance to win and is posted at 4 to 1 odds. (For every dollar Gavin bets, he receives \$5 if the horse wins and nothing if it loses.) He can either bet on this horse or keep his money in his pocket. Gavin decides that he has a square-root utility for money.
- (a) What fraction of his money should Gavin bet on No Arbitrage?  
(b) What is the implied winning payoff of a \$1 bet against No Arbitrage?
17. (General risk-neutral pricing) We can transform the log-optimal pricing formula into a risk-neutral pricing equation. From the log-optimal pricing equation we have

$$P = E\left(\frac{d}{R^*}\right)$$

where  $R^*$  is the return on the log-optimal portfolio. We can then define a new expectation operation  $\hat{E}$  by

$$\hat{E}(x) = E\left(\frac{Rx}{R^*}\right).$$

This can be regarded as the expectation of an artificial probability. Note that the usual rules of expectation hold. Namely:

- (a) If  $x$  is certain, then  $\hat{E}(x) = x$ . This is because  $E(1/R^*) = 1/R$ .

- (b) For any random variables  $x$  and  $y$ , there holds  $\hat{E}(ax + by) = a\hat{E}(x) + b\hat{E}(y)$ .
- (c) For any nonnegative random variable  $x$ , there holds  $\hat{E}(x) \geq 0$ .

Using this new expectation operation, with the implied artificial probabilities, show that the price of any security  $d$  is

$$P = \frac{\hat{E}(d)}{R}.$$

## References

The systematic use of expected utility as a basis for financial decision making was originated by von Neumann and Morgenstern in [1]. Another set of axioms is due to Savage [2]. The practical application of the theory was elaborated in [3]. For a comprehensive treatment explicitly aimed at finance problems, see [4]. The presentation of the second half of this chapter, related to linear pricing, draws heavily on the first chapter of [5]. The idea of linear pricing was developed in [6]. The use of the log-optimal portfolio to determine prices is explained in [7]. The idea of risk-neutral evaluation emerged from the pioneering approach to options by Black and Scholes in [8] and was formalized explicitly in [9]. The concept was generalized in [10] and now is a fundamental part of modern investment science. The concept of zero-level pricing was studied in [11–13]. For the notion of universality, see [14]. Also see [15].

1. von Neumann, J., and O. Morgenstern (1944), *Theory of Games and Economic Behavior*, Princeton University Press, Princeton, NJ.
2. Savage, L. J. (1954, 1972), *Foundations of Statistics*, Wiley, New York, 1954; 2nd ed., Dover, New York, 1972.
3. Luce, R. D., and H. Raiffa (1957), *Games and Decisions*, Wiley, New York.
4. Ingersoll, J. E., Jr. (1987), *Theory of Financial Decision Making*, Rowman and Littlefield, Savage, MD.
5. Duffie, D. (2001), *Dynamic Asset Pricing*, 3rd ed., Princeton University Press, Princeton, NJ.
6. Cox, J., S. Ross, and M. Rubinstein (1979), "Option Pricing: A Simplified Approach," *Journal of Financial Economics*, **7**, 229–263.
7. Long, J. B., Jr. (1990), "The Numeraire Portfolio," *Journal of Financial Economics*, **26**, 29–69.
8. Black, F., and M. Scholes (1973), "The Pricing of Options and Corporate Liabilities," *Journal of Political Economy*, **81**, 637–654.
9. Ross, S. (1961), "A Simple Approach to the Valuation of Risky Streams," *Journal of Business*, **34**, 411–433.
10. Harrison, J. M., and D. Kreps (1979), "Martingales and Arbitrage in Multiperiod Securities Markets," *Journal of Economic Theory*, **20**, 381–408.
11. Smith, J. E., and R. F. Nau (1995), "Valuing Risky Projects: Options Pricing Theory and Decision Analysis," *Management Science*, **41**, 795–816.
12. Holtan, H. M. (1997), *Asset Valuation and Optimal Portfolio Choice in Incomplete Markets*, Ph.D. dissertation, Department of Engineering-Economic Systems, Stanford University, Stanford, CA.
13. Luenberger, D. G. (1998), *Investment Science*, 1st ed., Oxford University Press, New York.
14. Luenberger, D. G. (2002), "Arbitrage and Universal Pricing," *Journal of Economic Dynamics and Control*, **26**, 1613–1628.
15. Guu, S. M., and J. N. Wang (2008), "Zero-Level Pricing and the HARA Utility Functions," *Journal of Optimization Theory and Applications*, **139**, 393–402.



## PART III

# DERIVATIVE SECURITIES





# 12

## FORWARDS, FUTURES, AND SWAPS

**A** derivative security is a security whose payoff is explicitly tied to the value of some other variable. In practice, however, this broad definition is often restricted to securities whose payoffs are explicitly tied to the price of some other financial security. A hypothetical example of such a derivative security is a certificate that can be redeemed in 6 months for an amount equal to the price, then, of a share of IBM stock. The certificate is a derivative security since its payoff depends on the future price of IBM. Most real derivatives are fashioned to have important risk control features, and the payoff relation is more subtle than that of the hypothetical certificate example. A more realistic example is a **forward contract** to purchase 2,000 pounds of sugar at 12 cents per pound in 6 weeks. There is no reference to a payoff—the contract just guarantees the purchase of sugar—but in fact a payoff is implied. The payoff is determined by the price of sugar in 6 weeks. If the price of sugar then were, say, 13 cents per pound, the contract would have a value of 1 cent per pound, or \$20, since the owner of the contract could buy sugar at 12 cents according to the contract and then turn around and sell that sugar in the sugar market at 13 cents. The contract is a derivative security because its value is derived from the price of sugar. Another realistic example is a contract that gives one the right, but not the obligation, to purchase 100 shares of GM stock for \$60 per share in exactly 3 months. This is an **option** to buy GM. The payoff of this option will be determined in 3 months by the price of GM stock at that time. If GM is selling then for \$70, the option will be worth \$1,000 because the owner of the option could at that time purchase 100 shares of GM for \$60 per share according to the option contract, and immediately sell those shares for \$70 each. As a final example of a derivative security, suppose you take

out a mortgage whose interest rate is adjusted periodically according to a weighted average of the rates on new mortgages offered by major banks. Your mortgage is a derivative security since its value at later times is determined by other financial prices, namely, prevailing interest rates.

As mentioned earlier, the payoff of a derivative security is usually based on the price of some other financial security. In the foregoing examples these were the price of IBM shares, the price of sugar, the price of GM shares, and the prevailing interest rates. The security that determines the value of a derivative security is called the **underlying security**. However, according to the broad definition, derivatives may have payoffs that are functions of nonfinancial variables, such as the weather or the outcome of an election. The main point is that the payments derived from a derivative security are deterministic functions of some other variable whose value will be revealed before or at the time of the payoff.

The main types of derivative securities are forward contracts, futures contracts, options, options on futures, and swaps.<sup>1</sup> Such securities play an important role in everyday commerce, since they provide effective tools for hedging risks involving the underlying variables. For example, a business that deals with a lot of sugar—perhaps a sugar producer, a processor, a marketeer, or a commercial user—typically faces substantial risks associated with possible sugar price fluctuations. Such users can control that risk through the use of derivative securities (in this case mainly through the use of sugar futures contracts). Indeed, the primary function of derivative securities in a portfolio—for businesses, institutions, or individuals—is to control risk.

This third part of the text addresses several aspects of derivative securities. First, these chapters explain what these different types of securities are; that is, how forwards, futures, swaps, and options are structured. Second, these chapters show, through theory and example, how derivative securities are used to control risk; that is, how derivatives can enhance the overall structure of a portfolio that contains risky components. Third, these chapters present the special pricing theory that applies to derivative securities. This is the aspect that receives the most attention in the text. Finally, an important technical subject presented in this part of the text is concerned with how to model security price fluctuations. This is the primary topic of the next chapter. This current chapter is devoted to forward and futures contracts, which are among the simplest and most useful derivative securities.

## 12.1 Pricing Principles

Rational derivative pricing is hailed as one of the great achievements of modern finance theory, and indeed this theory has led to the development of a huge industry

---

<sup>1</sup> In addition to the primary types listed here, there are *many* other derivative securities, such as variable-rate preferred stock, variable-rate mortgages, prime-rate loans, and LIBOR-based notes. New derivative securities are created and marketed every year by financial institutions. Fortunately most of these various financial products can be analyzed by using just a few common principles.

that every day directs the flow of billions or trillions of dollars. Perhaps it is not surprising that the underlying theory is regarded by many people as esoteric, something only for Ph.D.'s.

Actually, most derivative pricing theory is based on just three remarkably simple and intuitive principles, which, when used in combination, have great power.

These principles apply to well-functioning markets that satisfy a set of **perfect market** assumptions. One set of assumptions is that it is possible to buy, sell, or short-sell any asset; there are no transaction costs or taxes; no one person's action influences prices; every asset is infinitely devisable (that is, any portion of an asset is also an asset); the payoffs can be linearly combined (like cash); and no arbitrage opportunity exists in the market. Here are the principles:

- 1. Use the market.**
- 2. Discount certain cash at the current rate of interest.**
- 3. Use linear pricing.**

To define the first principle, suppose your rich uncle promises that he will give you an ounce of gold 1 year from now (and you know that he will keep his promise). Since the price of gold seems to fluctuate a lot, you would like to determine how valuable that volatile future gift is today. One way to arrive at a value is to estimate what the price of gold will be a year from now. Suppose it is \$1,000 today. You might reasonably conjecture that it may rise from this to perhaps \$1,100 or \$1,200 or even higher. Or, due to unforeseen circumstances, it may drop to \$900 or even \$875. You could talk to a number of experts to get their opinions. Eventually, you might settle on \$1,150 as a good estimate of the expected value. Then, because it is risky, you will want to discount the value heavily to get the current value of the gift. If you decide to use 15%, you will conclude that the value of this gift is  $\$1,150/1.15 = \$1,000$ , which happens (by chance?) to be today's gold price.

This forward estimation followed by discounting is a common procedure in asset pricing. However, there is a better way for derivatives.

We may conclude immediately that the promise of an ounce of gold in a year is worth exactly \$1,000—today's gold price. The reason is that your uncle can fund his promise by buying gold today for \$1,000 and delivering it to you next year.<sup>2</sup> Likewise you could today transform the gift into \$1,000 by shorting an ounce of gold and then clearing (that is, repaying) that short next year when your uncle gives you the ounce.

Using the market, we bypass the forecasting and discounting procedures. That is the first principle: **Use the market.**

For the next principle, suppose your rich uncle promises to give you \$1,000 in 1 year. What is the current value of this gift? This is much easier. The value is \$1,000 discounted by 1 year's interest. Thus, if the interest rate is 10%, the value is \$909. This is the second principle: **Discount certain cash at the rate of interest.** This principle can be regarded as a special case of the first principle if it is known

---

<sup>2</sup> We are assuming that there is no cost for storing gold.

that there is a risk-free asset in the market, for then an asset that pays \$1.00 is worth what a future dollar costs, namely, the discount factor.

Finally, the third principle is that linear pricing applies to combinations of assets. Specifically, if asset A has value  $V_A$  and asset B has value  $V_B$ , then the value of  $a$  units of A and  $b$  units of B is  $aV_A + bV_B$ . Hence, this is the third principle: **Use linear pricing.**

It may seem surprising that such simple and intuitive rules—market, discount, and linear—form the essential core of modern derivative pricing.<sup>3</sup> As we progress through the chapters of this part of the text, the simplicity of the underlying principles should be kept in mind. There are, however, cases when the market is not perfect (such as when shorting is not possible); we must then employ additional methods.

In Section 12.3 we shall apply the three principles to determine quickly the forward price of an asset.

## 12.2 Forward Contracts

Forward and futures contracts are closely related structures, but forward contracts are the simpler of the two. A **forward contract** on a commodity is a contract to purchase or sell a specific amount of the commodity at a specific price and at a specific time in the future. For example, a typical forward contract might be to purchase 100,000 pounds of sugar at 12 cents per pound on the 15th of March next year. The contract is between two parties, the buyer and the seller. The buyer is said to be **long** 100,000 pounds of sugar, and the seller is said to be **short**. Being long or short a given amount is the **position** of the party. Forward contracts for commodities have existed for thousands of years, for they are indeed a natural adjunct to commerce. Both suppliers and consumers of large quantities of a commodity frequently find it advantageous to lock in the price associated with a future commodity delivery.

A forward contract is specified by a legal document, the terms of which bind the two parties involved to a specific transaction in the future. However, a forward contract on a priced asset, such as sugar, is also a financial instrument, since it has an intrinsic value determined by the market for the underlying asset. Forward contracts have been extended in modern times to include underlying assets other than physical commodities. For example, many corporations use forward contracts on foreign currency or on interest rate instruments.

Most forward contracts specify that all claims are settled at the defined future date (or dates); both parties must carry out their side of the agreement at that time. Almost always, the initial payment associated with a forward contract is zero. Neither party pays any money to obtain the contract (although a security deposit is sometimes required of both parties). The **forward price** is the price that applies at delivery. This price is negotiated so that the initial payment is zero; that is, the *value* of the contract is zero when it is initiated.

---

<sup>3</sup> In fact, all three of these principles can be derived from the single principle termed the *law of one price*, which states that if two assets have the same payoffs, the prices must be equal. However, for most purposes it is easier to go directly to the three stated principles.

The open market for immediate delivery of the underlying asset is called the **spot market**. This is distinguished from the **forward market**, which trades contracts for future delivery. During the course of a forward contract, the spot market price may fluctuate. Hence, although the initial value of a forward contract is zero, its later values will vary as a function of the spot price of the underlying asset (or assets). Later we shall explore the relation between the current value and the forward price.

## Forward Interest Rates

We discussed a rather advanced form of forward contract in Chapter 4 when studying the term structure of interest rates. The forward rate was defined as the rate of interest associated with an agreement to loan money over a specified interval of time in the future. It may not be apparent how to arrange for such a loan using standard financial securities; but actually it is quite simple, as the following example illustrates.

**Example 12.1 (A T-bill forward)** Suppose that you wish to arrange to loan money for 6 months beginning 3 months from now. Suppose that the forward rate for that period is 10%. A suitable contract that implements this loan would be an agreement for a bank to deliver to you, 3 months from now, a 6-month Treasury bill (that is, a T-bill with 6 months to run from the delivery date). The price would be agreed upon today for this delivery, and the Treasury bill would pay its face value of, say, \$1,000 at maturity. The correct price for a Treasury bill of face value \$1,000 would be determined by the forward rate, which is 10% in annual terms, or 5% for 6 months. Hence the value of the T-bill would be  $\$1,000/1.05 = \$952.38$ , so this is the price that today you would agree to pay in 3 months when the T-bill is delivered to you. Six months later you receive the \$1,000 face value. Hence, overall, you have loaned \$952.38 for 6 months, with repayment of \$1,000. This agreement exactly parallels that of other forward contracts, the special feature being that the underlying asset to be delivered is a T-bill. The price associated with this contract directly reflects the forward interest rate.

The forward rates can be determined from the term structure of interest rates, which in turn can be determined from current bond prices. These forward rates are basic to the pricing of forward contracts on *all* commodities and assets because they provide a point of comparison. The payoff associated with a given forward contract on, say, sugar can be compared with one associated with pure lending and borrowing. Consistency (or lack of arbitrage opportunities) dictates the (theoretical) forward price, as we show next.

## 12.3 Forward Prices

As discussed earlier, there are two prices or values associated with a forward contract. The first is the **forward price**  $F$ . This is the delivery price of a unit of the underlying asset to be delivered at a specific future date. It is the delivery price that would be

specified in a forward contract written today. The second price or value of a forward contract is its current value, which is denoted by  $f$ . The forward price  $F$  is determined such that  $f = 0$  initially, so that no money need be exchanged when completing the contract agreement. After the initial time, the value  $f$  may vary, depending on variations of the spot price of the underlying asset, the prevailing interest rates, and other factors. Likewise the forward price  $F$  of new contracts with delivery terms identical to that of the original contract will also vary.

In this section we determine the theoretical forward price  $F$  associated with a forward contract written at time  $t = 0$  to deliver an asset at time  $T$ . Our analysis depends on the standard assumptions that there are no transaction costs, and that assets can be divided arbitrarily. Also we assume initially that it is possible to store the underlying asset without cost and that it is possible to sell the asset short. Later we will allow for storage costs, but still require that it be possible to store the underlying asset for the duration of the contract. This is a good assumption for many assets, such as gold or sugar or T-bills, but perhaps not good for perishable commodities such as oranges.

Suppose that at time  $t = 0$  the underlying asset has spot price  $S$  and a forward contract is being designed today for delivery at time  $T$ . How can we determine the value of the forward contract? We use the pricing principles. At time  $T$  we will acquire the commodity. By the first principle, the value at time 0 of this acquisition is  $S$ , the current spot price. This is how much we would pay at time 0 for delivery at  $T$ . However, the contract specifies that we pay at time  $T$ , not time 0. According to the second principle (and standard discounting theory) this delay in payment means that today the contract is worth  $S/d(0, T)$ , where  $d(0, T)$  is the discount factor between 0 and  $T$ . Hence,

$$F = S/d(0, T).$$

This result can also be obtained, more laboriously, by an argument based on the assumption that the price should not introduce an arbitrage opportunity. We present this argument in detail, not because it is really necessary here, but to set the stage for later situations where the present-value formula breaks down because of a market imperfection.

**Forward price formula** *Suppose an asset can be stored at zero cost and also sold short. Suppose the current spot price (at  $t = 0$ ) of the asset is  $S$ . The theoretical forward price  $F$  (for delivery at  $t = T$ ) is*

$$F = S/d(0, T) \tag{12.1}$$

where  $d(0, T)$  is the discount factor between 0 and  $T$ .

**Proof:** First suppose to the contrary that  $F > S/d(0, T)$ . Then we construct a portfolio as follows: At the present time borrow  $S$  amount of cash, buy one unit of the underlying asset on the spot market at price  $S$ , and take a one-unit short position in the forward market. The total cost of this portfolio is zero. At time  $T$  we deliver the asset (which we have stored), receiving a cash amount  $F$ , and we repay our loan in the amount  $S/d(0, T)$ . As a result we obtain a positive profit of

**TABLE 12.1**

<b>At <math>t = 0</math></b>	<b>Initial cost</b>	<b>Final receipt</b>
Borrow $\$S$	$-S$	$-S/d(0, T)$
Buy 1 unit and store	$S$	0
Short 1 forward	0	$F$
Total	0	$F - S/d(0, T)$

**TABLE 12.2**

<b>At <math>t = 0</math></b>	<b>Initial cost</b>	<b>Final receipt</b>
Short 1 unit	$-S$	0
Lend $\$S$	$S$	$S/d(0, T)$
Go long 1 forward	0	$-F$
Total	0	$S/d(0, T) - F$

$F - S/d(0, T)$  for zero net investment. This is an arbitrage, which we assume is impossible. The details of these transactions are shown in Table 12.1.

If  $F < S/d(0, T)$ , we can construct the reverse portfolio. However, this requires that we short one unit of the asset. The shorting is executed by borrowing the asset from someone who plans to store it during this period, then selling the borrowed asset at the spot price, and replacing the borrowed asset at time  $T$ . The arbitrage portfolio is constructed by shorting one unit, lending the proceeds  $S$  from time 0 to  $T$ , and taking a one-unit long position in the forward market. The net cash flow at time zero of this portfolio is zero. At time  $T$  we receive  $S/d(0, T)$  from our loan, pay  $F$  to obtain one unit of the asset, and we return this unit to the lender who made the short possible. The details are shown in Table 12.2.

Our profit is  $S/d(0, T) - F$  (which we might share with the asset lender for making the short possible).

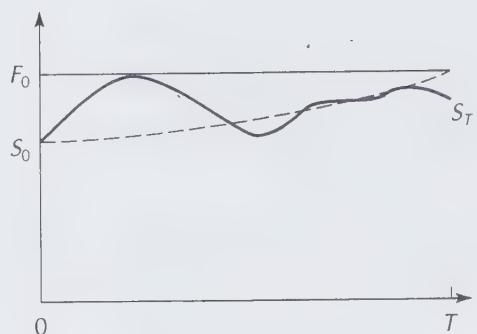
Since either inequality leads to an arbitrage opportunity, equality must hold. ■

The relationship between the spot price  $S$  and the forward price  $F$  is illustrated in Figure 12.1. The spot price starts at  $S(0)$  and varies randomly, arriving at  $S(T)$ . However, the forward price at time zero is based on extrapolating the current spot price forward at the prevailing rate of interest.

**Example 12.2 (Copper forward)** A manufacturer of heavy electrical equipment wishes to take the long side of a forward contract for delivery of copper in 9 months. The current price of copper is 84.85 cents per pound, and 9-month T-bills are selling at 970.87. What is the appropriate forward price of the copper contract?

If we ignore storage costs and use the T-bill rate, the appropriate price is  $84.85/970.87 = 87.40$  cents per pound.

**FIGURE 12.1 Forward price.** The forward price at time zero is equal to the projected future value of cash of amount  $S(0)$ .



**Example 12.3 (Continuous-time compounding)** If there is a constant interest rate  $r$  compounded continuously, the forward rate formula becomes

$$F = Se^{rT}.$$

The discount rate  $d(0, T)$  used in the forward price formula should be the one consistent with one's access to the interest rate market. Professional traders of forwards and futures commonly use the **repo rate** associated with repurchase agreements. (These are agreements to sell a security and repurchase it a short time later for a slightly higher price.) This repo rate is only slightly higher than the Treasury bill rate.

## Costs of Carry

The preceding analysis assumed that there are no storage costs associated with holding the underlying asset. This is not always the case. Holding a physical asset such as gold entails storage costs, such as vault rental and insurance fees. Holding a security may, alternatively, entail negative costs, representing dividend or coupon payments. These costs (or incomes) affect the theoretical forward price.

We shall use a discrete-time (multiperiod) model to describe this situation. The delivery date  $T$  is  $M$  periods (say, months) in the future. We assume that storage is paid periodically, and we measure time according to these periods. The carrying cost is  $c(k)$  per unit for holding the asset in the period from  $k$  to  $k + 1$  (payable at the beginning of the period). The forward price of the asset is then determined by the structure of the forward interest rates applied to the holding costs and the asset itself.

**Forward price formula with carrying costs** Suppose an asset has a holding cost of  $c(k)$  per unit in period  $k$ , and the asset can be sold short. Suppose the initial spot price is  $S$ . Then the theoretical forward price is

$$F = \frac{S}{d(0, M)} + \sum_{k=0}^{M-1} \frac{c(k)}{d(k, M)}, \quad (12.2)$$

**TABLE 12.3**  
**DETAILS OF ARBITRAGE**

Time 0 action	Time 0 cost	Time $k$ cost	Receipt at time $M$
Short 1 forward	0	0	$F$
Borrow $\$S$	$-S$	0	$\frac{-S(0,M)}{d(0,M)}$
Buy 1 unit spot	$S$	0	0
Borrow $c(k)$ 's forward	$-c(0)$	$-c(k)$	$-\sum_{k=0}^{M-1} \frac{c(k)}{d(k,M)}$
Pay storage	$c(0)$	$c(k)$	0
Total	0	0	$F - \frac{S}{d(0,M)} - \sum_{k=0}^{M-1} \frac{c(k)}{d(k,M)}$

where  $d(k, M)$  is the discount factor from  $k$  to  $M$ . Equivalently,

$$S = -\sum_{k=0}^{M-1} d(0,k)c(k) + d(0,M)F. \quad (12.3)$$

**Proof:** The simple version of the proof is this: Buy one unit of the commodity on the spot market and enter a forward contract to deliver one unit at time  $T$ . The cash flow stream associated with this is  $(-S - c(0), -c(1), -c(2), \dots, -c(M-1), F)$ . The present value of this stream must be zero, and this gives the stated formula for  $F$ . A detailed proof based on the no-arbitrage condition can also be constructed.

The details are shown in Table 12.3 for one direction. (See Exercise 5.) ■

The alternative formula (12.3) is obtained from formula (12.2) by multiplying through by  $d(0, M)$  and using the fact that  $d(0, M) = d(0, k)d(k, M)$  for any  $k$ . This alternative formula is probably the simplest to understand, since it is a standard present value equation. We recognize that we can buy the commodity at price  $S$  and deliver it according to a forward contract at time  $M$  in a completely deterministic fashion. The cash flow incurred while holding the commodity will be the carrying charges and the delivery price. The present value of this stream must equal the price  $S$ .

**Example 12.4 (Sugar with storage cost)** The current price of sugar is 12 cents per pound. We wish to find the forward price of sugar to be delivered in 5 months. The carrying cost of sugar is .1 cent per pound per month, to be paid at the beginning of the month, and the interest rate is constant at 9% per annum.

The interest rate is  $.09/12 = .0075$  per month. The reciprocal of the 1-month discount rate (for any month) is 1.0075. Therefore we find

$$\begin{aligned} F &= (1.0075)^5 (.12) + [(1.0075)^5 + (1.0075)^4 + (1.0075)^3 \\ &\quad + (1.0075)^2 + 1.0075](.001) \\ &= .1295 = 12.95 \text{ cents.} \end{aligned}$$

**Example 12.5 (A bond forward)** Consider a Treasury bond with a face value of \$10,000, a coupon of 8%, and several years to maturity. Currently this bond is selling for \$9,260, and the previous coupon has just been paid. What is the forward price for delivery of this bond in 1 year? Assume that interest rates for 1 year out are flat at 9%.

We recognize that there will be two coupons before delivery: one in 6 months and just prior to delivery. Hence using the present-value form (12.3) and a 6-month compounding convention, we have immediately

$$\$9,260 = \frac{F + \$400}{(1.045)^2} + \frac{\$400}{1.045}.$$

This can be solved [or turned around to the form (12.2)] to give

$$F = \$9,260(1.045)^2 - \$400 - \$400(1.045) = \$9,294.15$$

(in decimal form, not 32nd's).

## Tight Markets

At any one time it is possible to define several different forward contracts on a given commodity, each contract having a different delivery date. If the commodity is a physical commodity such as soybean meal, the preceding theory implies that the forward prices of these various contracts will increase smoothly as the delivery date is increased because the value of  $F$  in (12.2) increases with  $M$ . In fact, however, this is frequently *not* the case.

Consider, for example, the prices for soybean contracts shown in Table 12.4. This table<sup>4</sup> shows that the prices actually decrease with time over a certain range. How do we explain this? Certainly the holding cost for soybean meal is not negative. In fact, holders of soybean meal are giving up an opportunity to make arbitrage profit.

To verify this opportunity, note that someone, say, a farmer with soybean meal could sell it now (in December) at \$188.20 and arrange now to buy it back in March at \$184.00, thereby making a sure profit and avoiding any holding costs that would

**TABLE 12.4**  
**SOYBEAN MEAL FORWARD PRICES**

Dec	188.20	Aug	185.50
Jan	185.60	Sept	186.20
Mar	184.00	Oct	188.00
May	183.70	Dec	189.00
July	184.80		

*The delivery prices do not increase continuously as the delivery date is increased.*

<sup>4</sup> These are actually futures market prices, but they can be assumed to be forward prices.

otherwise be incurred. Why does the farmer not do this? The reason is that soybean meal is frequently in short supply; those that hold it do so because they need it to supply other contracts or for their own use. It is true that they could make a small profit by selling their holdings and purchasing a forward contract, but this small potential profit is less than the costs incurred by not having soybean meal on hand.

Likewise, arbitrageurs are unable to short the commodity contract because no one will lend them soybean meal. Hence the theoretical price relationship that assumes that shorting is possible does not apply.

The theoretical relation does hold in one direction as long as storage is possible. This is the case for most assets (including soybean meal). When storage is possible, the first direction of the proofs of equations (12.1) and (12.2) applies. In other words,

$$F \leq \frac{S}{d(0,M)} + \sum_{k=0}^{M-1} \frac{c(k)}{d(k,M)} \quad (12.4)$$

must hold if there are no arbitrage opportunities.

Shorting, on the other hand, relies on there being a positive amount of storage available for borrowing over the period from 0 to  $T$ . Someone, or some group, must plan on having excess stocks over this entire period, no matter how the market changes. If stocks are low, or potentially low, short selling at the spot price is essentially infeasible. That means that the second direction of the proofs of (12.1) and (12.2) does not apply. Hence only the inequality (12.4) can be inferred. As shown by the example of soybean meal, this is, in fact, a fairly common situation.

The inequality can be converted to an equality by the artifice of defining a **convenience yield**, which measures the benefit of holding the commodity. In the case of soybean meal, for example, the convenience yield may represent the value of having meal on hand to keep a farm operating. The convenience yield can be thought of as a negative holding cost, so if incorporated into equation (12.4), it reduces the right-hand side to the point of equality. One way to incorporate it is to modify inequality (12.4) as

$$F = \frac{S}{d(0,M)} + \sum_{k=0}^{M-1} \frac{c(k)}{d(k,M)} - \sum_{k=0}^{M-1} \frac{y}{d(k,M)},$$

where  $y$  is the convenience yield per period.

## Investment Assets

It is useful to distinguish between two types of assets: investment assets and consumption assets. Investment assets are held by a significant number of people for investment purposes. Gold is a good example, as are silver and various securities. On the other hand, consumption assets are held primarily for consumption. Eggs are a good example, as are soybean meal and oil. Consumption assets are often in tight supply and hence in some situations cannot be sold short. On the other hand, someone who holds an investment asset is likely to be willing to sell some of that asset if there is opportunity for guaranteed profit.

At $t = 0$	Initial cost	Final receipt
Borrow 1 unit	0	
Sell 1 unit	- $S$	0
Lend $\$S$	$S$	$S/d(0, T)$
Go Long 1 forward	0	- $F$
		Return 1 unit
Total	0	$S/d(0, T) - F$

(a)

At $t = 0$	Initial cost	Final receipt
Sell 1 unit	- $S$	0
Lend $\$S$	$S$	$S/d(0, T)$
Go Long 1 forward	0	- $F$
Total	0	$S/d(0, T) - F$

(b)

**FIGURE 12.2 Construction of an arbitrage when  $S/d(0, T) - F > 0$  for two cases.**  
(a) Short-selling of the asset is possible; (b) the asset is held as an investment.

Figure 12.2 shows the detail of transactions that produce a sure profit when a forward price satisfies  $F < S/d(0, T)$ . There are two cases: (a) when shorting is possible, and (b) when the asset  $S$  is an investment asset. Notice that the pattern of cost and receipt columns are identical. In case (b) however, it is possible to sell without borrowing the asset.

It follows that for an investment asset, the forward price will adjust to eliminate the arbitrage possibility, and hence  $F \geq S/d(0, T)$ . A similar analysis applies if there are carrying costs.

## 12.4 The Value of a Forward Contract

Suppose a forward contract was written in the past with a delivery price of  $F_0$ . At the present time  $t$  the forward price for the same delivery date is  $F_t$ . We would like to determine the current value  $f_t$  of the initial contract. This value is given by the following statement.

**The value of a forward** Suppose a forward contract for delivery at time  $T$  in the future has a delivery price  $F_0$  and a current forward price  $F_t$ . The value of the contract is

$$f_t = (F_t - F_0)d(t, T),$$

where  $d(t, T)$  is the risk-free discount factor over the period from  $t$  to  $T$ .

**Proof:** Consider forming the following portfolio at time  $t$ : one unit long of a new forward contract with delivery price  $F_t$  maturing at time  $T$ , and one unit short of

the old contract with delivery price  $F_0$ . The initial cash flow of this portfolio is  $f_t$ . The final cash flow at time  $T$  is  $F_0 - F_t$ . This is a completely deterministic stream, because the short and long delivery requirements cancel. The present value of this portfolio is  $f_t + (F_0 - F_t)d(t, T)$ , and this must be zero. The stated result follows immediately. ■

## 12.5 Swaps\*

Motivating most investment problems is a desire to transform one cash flow stream into another by appropriate market or technological activity. A **swap** accomplishes this directly—for a swap is an agreement to exchange one cash flow stream for another. The attraction of this direct approach is evidenced by the fact that the swap market amounts to hundreds of billions of dollars. Swaps are often tailored for a specific situation, but the most common is the **plain vanilla swap**, in which one party swaps a series of variable payments for a series of fixed-level payments. It is this form that we consider in this section. As we shall see, such swaps can be regarded as a series of forward contracts, and hence they can be priced using the concepts of forwards.

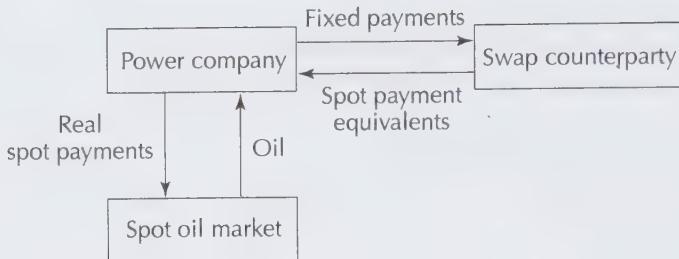
As an example, consider a plain vanilla interest rate swap. Party A agrees to make a series of semiannual payments to party B equal to a fixed rate of interest on a notional principal. (The term **notional principal** is used because there is no loan. This principal simply sets the level of the payments.) In return, party B makes a series of semiannual payments to party A based on a floating rate of interest (such as the current 6-month LIBOR rate) and the same notional principal. Usually, swaps are **netted** in the sense that only the difference of required payments is made by the party that owes the difference.

This swap might be motivated by the fact that party B has loaned money to a third party C under floating rate terms; but party B would rather have fixed payments. The swap with party A effectively transforms the floating rate stream to one with fixed payments.

As an example of a commodity swap, consider an electric power company that must purchase oil every month for its power generation facility. If it purchases oil on the spot market, the company will experience randomly fluctuating cash flows caused by fluctuating spot prices. The company may wish to swap this payment stream for one that is constant. It can do this if it can find a counterparty willing to swap. This is shown in Figure 12.3. The swap counterparty agrees to pay the power company the spot price of oil times a fixed number of barrels, and in return the power company pays a fixed price per barrel for the same number of barrels over the life of the swap. The variable cash flow stream is thereby transformed to a fixed stream.

### Value of a Commodity Swap

Consider an agreement where party A receives spot price for  $N$  units of a commodity each period while paying a fixed amount  $X$  per unit for  $N$  units. If the agreement is made for  $M$  periods, the net cash flow stream received by A is  $(S_1 - X, S_2 - X, \dots, S_M - X)$ .



**FIGURE 12.3 Commodity swap.** The power company buys oil on the spot market every month. The company arranges a swap with a counterparty (or a swap dealer) to exchange fixed payments for spot price payments. The net effect is that the power company has eliminated the variability of its payments.

$S_3 - X, \dots, S_M - X$ ) multiplied by the number of units  $N$ , where  $S_i$  denotes the spot price of the commodity at time  $i$ .

We can value this stream using the concepts of forward markets. At time zero the forward price of one unit of the commodity to be received at time  $i$  is  $F_i$ . This means that we are indifferent between receiving  $S_i$  (which is currently uncertain) at  $i$  and receiving  $F_i$  at  $i$ .  $F_i$  is, in fact, the market price of future delivery of  $S_i$ . By discounting back to time 0 we conclude that the current value of receiving  $S_i$  at time  $i$  is  $d(0, i)F_i$ , where  $d(0, i)$  is the discount factor at time 0 for cash received at  $i$ .

If we apply this argument each period, we find that the total value of the stream is

$$V = \sum_{i=1}^M d(0, i)(F_i - X)N. \quad (12.5)$$

Hence the value of the swap can be determined from the series of forward prices. Usually  $X$  is chosen to make the value zero, so that the swap represents an equal exchange.

**Example 12.6 (A gold swap)** Consider an agreement by an electronics firm to receive spot value for gold in return for fixed payments. We assume that gold is in ample supply and can be stored without cost—which implies that the swap formula takes an almost trivial form. In that case we know that the forward price is  $F_i = S_0/d(0, i)$ . Therefore, equation (12.5) becomes

$$V = \left[ MS_0 - \sum_{i=1}^M d(0, i)X \right] N.$$

The summation is identical to the value of the coupon payment stream of a bond. Using this fact, it is easy to convert the value formula to

$$V = \left\{ MS_0 - \frac{X}{C} [B(M, C) - 100d(0, M)] \right\} N, \quad (12.6)$$

where  $B(M, C)$  denotes the price (relative to 100) of a bond of maturity  $M$  and coupon  $C$  per period. Any value of  $C$  can be used. (See Exercise 8.)

## Value of an Interest Rate Swap

Consider a plain vanilla interest rate swap in which party A agrees to make payments of a fixed rate  $r$  of interest on a notional principal  $N$  while receiving floating-rate payments on the same notional principal for  $M$  periods. The cash flow stream received by A is  $(c_1 - r, c_2 - r, \dots, c_M - r)$  times the principal  $N$ . The  $c_i$ 's are the floating rates.

We can value the floating portion of this swap with a special trick derived from our knowledge of floating-rate bonds. (For a direct proof using forward pricing concepts, see Exercise 12.) The floating-rate cash flow stream is exactly the same as that generated by a floating-rate bond of principal  $N$  and maturity  $M$ , except that no final principal payment is made. We know that the initial value of a floating-rate bond (including the final principal payment) is par; hence the value of the floating-rate portion of the swap is par minus the present value of the principal received at  $M$ . In other words, the value of the floating-rate portion of the swap stream is  $N - d(0, M)N$ .

The value of the fixed-rate portion of the stream is the sum of the discounted fixed payments, discounted according to the current term structure discount rates. Hence overall, the value of the swap is<sup>5</sup>

$$V = \left[ 1 - d(0, M) - r \sum_{i=1}^M d(0, i) \right] N.$$

The summation can be reduced using the method in the gold swap example.

## 12.6 Basics of Futures Contracts

Because forward trading is so useful, it became desirable long ago to standardize the contracts and trade them on an organized exchange. An exchange helps define universal prices and provides convenience and security because individuals do not themselves need to find an appropriate counterparty and need not face the risk of counterparty default. Individual contracts are made with the exchange, the exchange itself being the counterparty for both long and short traders. But standardization presents an interesting challenge. Consider the likely mechanics of forward contract trading on an exchange. It is a relatively simple matter to standardize a set of delivery dates, quantities to be delivered, quality of delivered goods, and delivery locations (although there are some subtleties even in these items). But standardization of forward prices is impossible. To appreciate the issue, suppose that contracts were issued today at a delivery price of  $F_0$ . The exchange would keep track of all such contracts. Then

---

<sup>5</sup> Typically, account must be made for other details. For example, interest rates for fixed payments are usually quoted on the basis of 365 days per year, whereas for floating rates they are quoted on the basis of 360 days per year.

tomorrow the forward price might change and contracts initiated that day would have a different delivery price  $F_1$ . In fact, the appropriate delivery price might change continuously throughout the day. The thousands of outstanding forward contracts could each have a different delivery price, even though all other terms were identical. This would be a bookkeeping nightmare.

The way that this has been solved is through the brilliant invention of a **futures market** as an alternative to a forward market. Multiple delivery prices are eliminated by revising contracts as the price environment changes. Consider again the situation where contracts are initially written at  $F_0$  and then the next day the price for new contracts is  $F_1$ . At the second day, the clearinghouse associated with the exchange revises all the earlier contracts to the new delivery price  $F_1$ . To do this, the contract holders either pay or receive the difference in the two prices, depending on whether the change in price reflects a loss or a gain. Specifically, suppose  $F_1 > F_0$  and I hold a one-unit long position with price  $F_0$ . My contract price is then changed to  $F_1$  and I receive  $F_1 - F_0$  from the clearinghouse because I will later have to pay  $F_1$  rather than  $F_0$  when I receive delivery of the commodity.

The process of adjusting the contract is called **marking to market**. In more detail it works like this: An individual is required to open a **margin account** with a broker. This account must contain a specified amount of cash for each futures contract (usually on the order of 5–10% of the value of the contract). All contract holders, whether short or long, must have such an account. These accounts are marked to market at the end of each trading day. If the price of the futures contract (the price determined on the exchange) increased that day, then the long parties receive a profit equal to the price change times the contract quantity. This profit is deposited in their margin accounts. The short parties lose the same amount, and this amount is deducted from their margin accounts. Hence each margin account value fluctuates from day to day according to the change in the futures price. With this procedure, every long futures contract holder has the same contract, as does every short contract holder. At the delivery date, delivery is made at the futures contract price at that time, which may be quite different from the futures price at the time the contract was first purchased.

Actually, delivery of commodities under the terms of a futures contract is quite rare; over 90% of all parties close out their positions before the delivery date. Even commercial organizations that need the commodity for production frequently close out their long positions and purchase the commodity from their conventional suppliers on the spot market.

Futures prices are listed in financial newspapers such as *Investor's Business Daily*. An example listing for wheat futures is shown in Figure 12.4. The heading explains that a standard contract for wheat is for 5,000 bushels, and that prices are quoted in cents per bushel. Notice that the July 2011 price is lower than the May 2011 contract, indicating a tight market effect. Open interest is the total number of contracts outstanding. Delivery of the commodity must be made on a specific day of the delivery month.

Margin accounts not only serve as accounts to collect or pay out daily profits, they also guarantee that contract holders will not default on their obligations. Margin accounts usually do not pay interest, so the cash in these accounts is, in effect, losing money. However, many brokers allow Treasury bills or other securities, as well as

Contract		Open					
High	Low	Interest	Open	High	Low	Settle	Chg.

For Wednesday, May 18, 2011

### Grains

WHEAT	(CBOT)	-	5,000	bu	minimum	cents	per	bushel
950.75	539.75	Jul 11	219,010	816.75	822.25	767.00	817.00	+53.0
971.50	559.00	Sep 11	71,717	863.75	864.50	810.00	858.50	+51.3
986.75	342.00	Dec 11	94,399	908.00	908.50	858.25	902.75	+44.5
994.75	597.00	Mar 12	19,886	935.50	936.50	895.00	929.75	+35.0
993.25	604.00	May 12	6,769	944.00	944.00	908.25	940.00	+34.3
946.50	611.00	Jul 12	26,792	939.75	940.25	905.00	935.25	+34.3
961.00	644.25	Dec 12	11,795	957.00	960.00	930.00	957.50	+27.8
Est. Vol. 326,286		Vol. 108,137		open int 454,539 - 3,777				

**FIGURE 12.4 Wheat futures quotations for May 18, 2011.** The heading tells us that the commodity is wheat traded at the Chicago Board of Trade; that a contract is 5,000 bushels; and that prices are cents per bushel. The first two columns give the high and low prices over the lifetime of the contract. The next column specifies the delivery month. The second-to-last column gives the settled price for the day.

(Source: *Investor's Business Daily*, May 19, 2011.)

cash, to serve as margin, so interest can be earned indirectly. If the value of a margin account should drop below a defined maintenance margin level (usually about 75% of the initial margin requirement), a **margin call** is issued to the contract holder, demanding additional margin. Otherwise the futures position will be closed out by taking an equal and opposite position.

**Example 12.7 (Margin)** Suppose that Mr. Smith takes a long position of one contract in corn (5,000 bushels) for March delivery at a price of \$2.10 (per bushel). And suppose the broker requires margin of \$800 with a maintenance margin of \$600.

The next day the price of this contract drops to \$2.07. This represents a loss of  $0.03 \times 5,000 = \$150$ . The broker will take this amount from the margin account, leaving a balance of \$650. The following day the price drops again to \$2.05. This represents an additional loss of \$100, which is again deducted from the margin account. At this point the margin account is \$550, which is below the maintenance level. The broker calls Mr. Smith and tells him that he must deposit at least \$50 in his margin account, or his position will be closed out, meaning that Mr. Smith will be forced to give up his contract, leaving him with \$550 in his account.

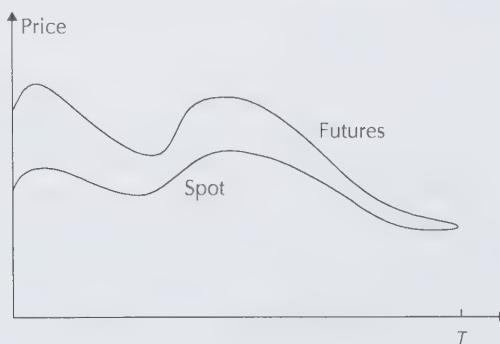
## 12.7 Futures Prices

There is, at any one time, only one price associated with a futures contract—the delivery price. The value of existing contracts is always zero because they are marked to market. The delivery price will in general be different from the spot price of the underlying asset, but the two must bear some relation to each other. In fact, as the maturity date approaches, the futures price and the spot price must approach each other, actually converging to the same value. This effect, termed **convergence**, is illustrated in Figure 12.5.

As a general rule we expect that the (theoretical) futures price should have a close relation to the forward price, the delivery price at which forward contracts would be written. Both are prices for future delivery. However, even if we idealize the mechanics of forward and futures trading by assuming no transactions costs and by assuming that no margin is required (or that margin earns competitive interest), there remains a fundamental difference between the cash flow processes associated with forwards and futures. With forwards, there is no cash flow until the final period, where either delivery is made or the contract is settled in cash according to the difference between the spot price and the previously established delivery price. With futures, there is cash flow every period after the first, the cash flow being derived from the most recent change in futures price. It seems likely that this difference in cash flow pattern will cause forward and futures prices to differ. In fact, however, under the assumption that interest rates are deterministic and follow expectations dynamics, as described in Chapter 4, the forward and futures prices must be identical if arbitrage opportunities are precluded. This important result is established here:

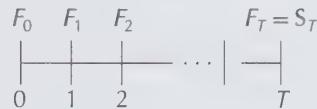
**Futures-forward equivalence** *Suppose that interest rates are known to follow expectations dynamics. Then the theoretical futures and forward prices of corresponding contracts are identical.*

**Proof:** Let  $F_0$  be the initial futures price (but remember that no payment is made initially). Let  $G_0$  be the corresponding forward price (to be paid at delivery time).



**FIGURE 12.5 Convergence of spot and futures prices.** The futures price converges to the spot price as time approaches the delivery date.

Assume that there are  $T + 1$  time points and corresponding futures prices, as indicated:



Let  $d(j, k)$  denote the discount rate at time  $j$  for a bond of unit face value maturing at time  $k$  (with  $j < k$ ).

We now consider two strategies for participation in the futures and forward markets, respectively.

### Strategy A

- At time 0: Go long  $d(1, T)$  futures.
- At time 1: Increase position to  $d(2, T)$
- ⋮
- At time  $k$ : Increase position to  $d(k+1, T)$
- At time  $T - 1$ : Increase position to 1.

The profit at time  $k + 1$  from the previous period is

$$(F_{k+1} - F_k)d(k+1, T).$$

As part of strategy A we invest this profit at time  $k + 1$  in the interest rate market until time  $T$ . It is thereby transformed to the final amount

$$\frac{d(k+1, T)}{d(k+1, T)}(F_{k+1} - F_k) = F_{k+1} - F_k.$$

The total profit from strategy A is therefore

$$\text{profit}_A = \sum_{k=0}^{T-1} (F_{k+1} - F_k) = F_T - F_0 = S_T - F_0.$$

Note that at each step before the end, there is zero net cash flow because all profits (or losses) are absorbed in the interest rate market. Hence a zero investment produces  $\text{profit}_A$ .

**Strategy B** Take a long position in one forward contract. This requires no initial investment and produces a profit of

$$\text{profit}_B = S_T - F_0.$$

We can now form a new strategy, which is  $A - B$ . This combined strategy also requires no cash flow until the final period, at which point it produces profit of  $G_0 - F_0$ . This is a deterministic amount, and hence must be zero if there is no opportunity for arbitrage. Hence  $G_0 = F_0$ . ■

When interest rates are not deterministic, the equivalence may not hold, but the equivalence is considered quite accurate for purposes of routine analysis. The result is important because it at least partially justifies simplifying an analysis of futures hedging by considering the corresponding hedge with forward contracts, where the cash flow occurs only at the delivery or settlement date.

**Example 12.8 (Wheat contracts)** In January a large producer of commercial flour and bread wishes to lock in the price for a large order of wheat. The producer would like to buy 500,000 bushels of wheat forward for May delivery. Although this producer could probably arrange a special forward contract, he decides instead to use the futures market, since it is organized and more convenient. The producer recognizes (and verifies) that the futures price is equal to the forward price he could negotiate.

The current futures (or forward) price for May delivery is \$3.30 per bushel. The size of a standard wheat futures contract is 5,000 bushels. Hence the producer decides that he needs 100 contracts.

Details of the futures market transaction are shown in Table 12.5. For simplicity this table shows accounting on a monthly basis, rather than on a daily basis.

The left part of the table shows the dates and the corresponding hypothetical prices (in cents) for a futures contract for May delivery. The next section, headed "Forward," shows the result of entering a forward contract for the delivery of 500,000 bushels of wheat in May, followed by the subsequent closing out of that contract so that delivery is not actually taken. There is no cash flow associated with this contract until May. Then there is the profit in May of 22 cents per bushel, or a total of \$110,000.

The next section of the table, headed "Futures contracts 1," shows the accounting details of entering a 100 contract long futures position in January and closing out this position in May. It is assumed that an account is established to hold all profits and losses. It is also assumed that the prevailing interest rate is 12%, or 1% per month, and that there are no margin requirements. Note that no money is required when the order is placed. A profit of \$50,000 is obtained in the second month because the futures price increased by 10 cents. This profit enters the account. The next month's

**TABLE 12.5**  
**FUTURES AND FORWARD TRANSACTIONS**

		<b>Forward</b>	<b>Futures contracts 1</b>				<b>Futures contracts 2</b>				
<b>Date</b>	<b>Price</b>		<b>Profit</b>	<b>Pos.</b>	<b>Profit</b>	<b>Interest</b>	<b>Balance</b>	<b>Pos.</b>	<b>Profit</b>	<b>Interest</b>	<b>Balance</b>
Jan 1	330		\$0	100	\$0	\$0	\$0	97	\$0	\$0	\$0
Feb 1	340		0	100	50,000	0	50,000	98	48,500	0	48,500
Mar 1	355		0	100	75,000	500	125,500	99	73,500	485	122,485
Apr 1	345		0	100	-50,000	1,225	76,755	100	-49,500	1,225	74,210
May 1	352	110,000	0	35,000	768	112,523		0	35,000	742	109,952
Total			\$110,000		\$110,000				\$107,500		

*The details of a forward contract, a fixed futures contract, and a futures contract strategy designed to mimic a forward are shown.*

balance reflects the additional profit and interest of the account. The total cash flow is \$110,000, exactly as in the case of the forward contract. However, because the cash flow occurs at various times, the actual final balance is \$112,523. (The result is more favorable in this case because prices rose early, but that is not the point.)

The third section of the table, headed “Futures contracts 2,” shows how futures can be used to duplicate a forward contract more precisely, by using the construction in the proof of the futures–forward equivalence result. Since interest is 1% a month, the discount rate increases by about 1% per month as well. Hence in this approach the producer initially goes long 97 contracts and increases this by 1 contract per month, finally reaching 100 contracts. Exactly the same accounting system is used as in the previous method. In this case the resulting final balance is \$109,952, which is very close to the \$110,000 figure obtained by a pure forward contract—the slight difference being due to rounding of the discount rate to even percentages so that integral numbers of contracts could be used.

This example illustrates that there is indeed a slight difference between forward and futures contract implementation if a constant contract level is used. In practice, however, the difference between using forward and futures contracts is small over short intervals of time, such as a few months. Furthermore, if interest rates are deterministic and follow expectations dynamics, then the difference between using futures and using forwards can be reduced to zero within rounding errors caused by the restriction to integral numbers of contracts.

## 12.8 Relation to Expected Spot Price\*

At time zero it is logical to form an opinion, or expectation, about the spot price of a commodity at time  $T$ . Is the current futures price for delivery at time  $T$  a good estimate of the future spot price; that is, is  $F = E(S_T)$ ?

If there were inequality, say,  $F < E(S_T)$ , a speculator might take a long position in futures and then at time  $T$  purchase the commodity at  $F$  according to the contract and sell the commodity at  $S_T$  for an expected profit of  $E(S_T) - F$ . If the inequality were in the other direction, the investor could carry out the reverse plan by taking a short position in futures. Hence speculators are likely to respond to any inequality.

Hedgers, on the other hand, participate in futures mainly to reduce the risks of commercial operations, not to speculate on commodity prices. Hence hedgers are unlikely to be influenced by small discrepancies between futures prices and expected spot prices.

Now suppose that there happen to be many more hedgers that are short in futures than those that are long. For the market to balance, speculators must enter the market and take long positions. They will do so only if they believe  $F < E(S_T)$ . Conversely, if there are more hedgers that are long in futures than those that are short, speculators will take the corresponding short position only if they believe  $F > E(S_T)$ .

The two situations have been given special names. If the futures price is below the expected future spot price, that is **normal backwardation**. If the futures price is above the expected future spot price, that is **contango**.

## 12.9 The Perfect Hedge

The primary use of futures contracts is to hedge against risk. Hedging strategies can be simple or complex; we shall illustrate some of the main approaches to their design in the remainder of this chapter.

The simplest hedging strategy is the **perfect hedge**, where the risk associated with a future commitment to deliver or receive an asset is completely eliminated by taking an **equal and opposite** position in the futures market. Equivalently, the hedge is constructed to effectively make anticipated future market purchases or sales immediately. This locks in the price of the futures transaction; there is absolutely no price risk. Such a strategy is possible only if there is a futures contract that exactly matches, with respect to the nature of the asset and the terms of delivery, the obligation that is being hedged.

**Example 12.9 (A wheat hedge)** Consider again the producer of flour and bread of Example 12.8. The producer has received a large order for delivery on May 20 at a specified price. To satisfy this order, the producer will purchase 500,000 bushels of wheat on the spot market shortly before the order is due. The producer has calculated its profit on the basis of current prices for wheat, but if the wheat price should measurably increase, the order may become unprofitable. The producer can hedge by taking an equal and opposite position in wheat futures. (That is, the producer is obligated to supply processed wheat, so it goes opposite the obligation and purchases wheat. Alternatively, the producer may think of it as *purchasing early* wheat that it must ultimately purchase.)

If we ignore the slight discrepancy between futures and forwards due to differences in cash flow timing, we can treat the futures contract just like a forward. The producer will close out the position in the futures market and then purchase wheat in the spot market. Since the price in the spot market will be the same as the closing futures price, the net effect is that the producer pays the original price of \$3.30 per bushel.

**Example 12.10 (A foreign currency hedge)** A U.S. electronics firm has received an order to sell equipment to a German customer in 90 days. The price of the order is specified as 500,000 euros, which will be paid upon delivery. The U.S. firm faces risk associated with the exchange rate between euros and U.S. dollars.

The firm can hedge this foreign exchange risk with four euros contracts (125,000 euros per contract) with a 90-day maturity date. Since the firm will be receiving euros in 90 days, it hedges by taking an equal and opposite position now—that is, it goes *short* four contracts. (Viewed alternatively, after receiving euros, the firm will want to sell them, so it sells them early by going short.)

## 12.10 The Minimum-Variance Hedge

It is not always possible to form a perfect hedge with futures (or forward) contracts. There may be no contract involving the exact asset whose value must be hedged, the

delivery dates of the available contracts may not match the asset obligation date, the amount of the asset obligated may not be an integral multiple of the contract size, there may be a lack of liquidity in the futures market, or the delivery terms may not coincide with those of the obligation. In these situations, the original risk cannot be eliminated completely with a futures contract, but usually the risk can be reduced.

One measure of the lack of hedging perfection is the **basis**, defined as the mismatch between the spot and futures prices. Specifically,

$$\text{basis} = \text{spot price of asset to be hedged} - \text{futures price of contract used.}$$

If the asset to be hedged is identical to that of the futures contract, then the basis will be zero at the delivery date. However, in general, for the reasons mentioned, the final basis may not be zero as anticipated. Usually the final basis is a random quantity, and this precludes the possibility of a perfect hedge. The basis risk calls for alternative hedging techniques.

One common method of hedging in the presence of basis risk is the minimum-variance hedge. The general formula for this hedge can be deduced quite readily. Suppose that at time zero the situation to be hedged is described by a cash flow  $x$  to occur at time  $T$ . For example, if the obligation is to purchase  $W$  units of an asset at time  $T$ , we have  $x = -WS$ , where  $S$  is the spot price of the asset at  $T$ . Let  $F$  denote the futures price of the contract that is used as a hedge, and let  $h$  denote the futures position taken. We neglect interest payments on margin accounts by assuming that all profits (or losses) in the futures account are settled at  $T$ . The cash flow at time  $T$  is therefore equal to the original obligation plus the profit in the futures account. Hence,

$$\text{cash flow} = y = x + (F_T - F_0)h.$$

We find the variance of the cash flow as

$$\text{var}(y) = E[x - \bar{x} + (F_T - \bar{F}_T)h]^2 = \text{var}(x) + 2\text{cov}(x, F_T)h + \text{var}(F_T)h^2.$$

This is minimized by setting the derivative with respect to  $h$  equal to zero. This leads to the following result:

**Minimum-variance hedging formula** *The minimum-variance hedge and the resulting variance are*

$$h = -\frac{\text{cov}(x, F_T)}{\text{var}(F_T)} \quad (12.7)$$

$$\text{var}(y) = \text{var}(x) - \frac{\text{cov}(x, F_T)^2}{\text{var}(F_T)}. \quad (12.8)$$

**Proof:** Setting to zero the derivative of the expression for  $\text{var}(y)$  with respect to  $h$ , we find  $2 \text{cov}(x, F_T) + 2 \text{var}(F_T)h = 0$ . This leads to the expression for  $h$ . Substitution of this  $h$  into the expression for  $\text{var}(y)$  gives

$$\begin{aligned}\text{var}(y) &= \text{var}(x) - 2 \frac{\text{cov}(x, F_T)^2}{\text{var}(F_T)} + \frac{\text{cov}(x, F_T)^2}{\text{var}(F_T)} \\ &= \text{var}(x) - \frac{\text{cov}(x, F_T)^2}{\text{var}(F_T)}.\end{aligned}$$

When the obligation has the form of selling a fixed amount  $W$  (that is,  $x = WS_T$ ) of an asset whose spot price is  $S_T$ , equation (12.7) becomes

$$h = -\beta W, \quad (12.9)$$

where

$$\beta = \frac{\text{cov}(S_T, F_T)}{\text{var}(F_T)}.$$

This, of course, reminds us of the general mean–variance formulas of Chapter 7; and indeed it is closely related to them.

**Example 12.11 (The perfect hedge)** As a special case, suppose that the futures commodity is identical to the spot commodity being hedged. In that case  $F_T = S_T$ . Suppose that the obligation is to sell  $W$  units of the commodity, so  $x = WS_T$ . In that case,  $\text{cov}(x, F_T) = \text{cov}(S_T, F_T)W = \text{var}(F_T)W$ . Therefore, according to equation (12.7), we have  $h = -W$ , and according to (12.8), we find  $\text{var}(y) = 0$ . In other words, the minimum-variance hedge reduces to the perfect hedge when the futures price is perfectly correlated with the spot price of the commodity being hedged.

**Example 12.12 (Hedging foreign currency with alternate futures)** The BIG H Corporation (a U.S. corporation) has obtained a large order from a small developing country, whose currency is the zee. Payment will be in 60 days in the amount of 1 million zees. BIG H would like to hedge the exchange risk, but, as is typical of developing countries, there is no forward contract for its currency. The vice president for finance of BIG H decides that the company can use a **cross hedge** by using Japanese yen, for, although the yen and the zee do not follow each other exactly, there is a fairly close relation—certainly closer than the zee and the dollar.

He notes that the current exchange rates are  $Z = .164$  dollar/zee and  $Y = .0125$  dollar/yen. Hence the exchange rate between yen and zees is  $Z/Y = .164/.0125 = 13.12$  yen/zee. Therefore 1 million zees is equivalent to 13.12 million yen at the current exchange rate. He deduces that an equal and opposite hedge would be to short 13.12 million yen.

An intern working at BIG H suggests that a minimum-variance hedge be considered as an alternative. The intern is given a few days to work out the details. He does some quick historical studies and estimates that the monthly fluctuations in the U.S. exchange rates for  $Z$  and  $Y$  are correlated with a correlation coefficient of

about .8. The standard deviation of these fluctuations is found to be about 3% of its value per month for yen and slightly less, 2.5%, for zees. In this problem the  $x$  of equation (12.7) denotes the dollar value of 1 million zees in 60 days, and  $F_T$  is the dollar value of a yen at that time. We may put  $x = Z \times 1$  million. The intern therefore estimates beta as

$$\beta = \frac{\text{cov}(Z, Y)}{\text{var}(Y)} = \frac{\sigma_{ZY}}{\sigma_Z \sigma_Y} \times \frac{\sigma_Z}{\sigma_Y} = \rho \frac{\sigma_Z}{\sigma_Y} = .8 \times \frac{.025Z}{.030Y}.$$

Hence the minimum-variance hedge is

$$\begin{aligned} h &= -\frac{\text{cov}(x, F)}{\text{var}(F)} = -\frac{\text{cov}(Z, Y) \times 1,000,000}{\text{var}(F)} \\ &= \left[ -.8 \times \frac{2.5}{3.0} \times 13.12 \times 1,000,000 \right] = -8.75 \text{ million yen.} \end{aligned}$$

The minimum-variance hedge is smaller than the equal and opposite hedge based on the exchange rates; it is reduced by the correlation coefficient and by the ratio of standard deviations.

We can go a bit further and find out how effective the hedge really is as compared to doing nothing. We have  $x = Z \times 1$  million. Hence  $\text{cov}(x, Y) = 1 \text{ million} \times \sigma_{ZY}$  and  $\sigma_x = 1 \text{ million} \times \sigma_Z$ . Combining these two, we have  $\text{cov}(x, Y) = \sigma_{ZY}\sigma_x/\sigma_Z$ . Using the minimum-variance hedging formula, we find

$$\begin{aligned} \text{var}(y) &= \text{var}(x) - \frac{\text{cov}(x, Y)^2}{\sigma_y^2} = \left[ 1 - \left( \frac{\sigma_{ZY}}{\sigma_Y \sigma_Z} \right)^2 \right] \text{var}(x) \\ &= \sqrt{1 - \rho^2} \text{var}(x). \end{aligned}$$

Thus,

$$\text{stdev}(y) = \left( \sqrt{1 - .8^2} \right) \text{stdev}(x) = .6 \times \text{stdev}(x).$$

Hence the minimum-variance hedge reduces risk by a factor of .6. A hedge with a lower risk would be obtained if a hedging instrument could be found that was more highly correlated with the zee.

**Example 12.13 (Changing portfolio beta with stock index futures)** Mrs. Smith owns a large portfolio that is heavily weighted toward high technology stocks. She believes that these securities will perform exceedingly well compared to the market as a whole over the next several months. However, Mrs. Smith realizes that her portfolio, which has a beta (with respect to the market) of 1.4, is exposed to a significant degree of market risk. If the general market declines, her portfolio will also decline, even if her securities do achieve significant excess return above that predicted by, say, CAPM, as she believes they will.

Mrs. Smith decides to hedge against this market risk. She can change the beta of her portfolio by selling some stock index futures. She might decide to construct a minimum-variance hedge of her \$2 million portfolio by shorting \$2 million  $\times 1.4 = \$2.8$  million of S&P 500 stock index futures with maturity in 120 days. Since the

normal beta of her portfolio is based on the S&P 500, this beta is the same beta as in the general equation, (12.9). The overall new beta of her hedged portfolio, after taking the short position in the stock index futures, is zero.

## 12.11 Optimal Hedging\*

Although the minimum-variance hedge is useful and fairly simple, it can be improved by viewing the hedging problem from a portfolio perspective. Suppose again that there is an existing cash flow commitment  $x$  at time  $T$ . And suppose that this will be hedged by futures contracts in the amount  $h$ , leading to a final cash flow of  $x + h(F_T - F_0)$ . If a utility function is assigned, it is appropriate to solve the problem<sup>6</sup>

$$\underset{h}{\text{maximize}} \quad E\{x + h(F_T - F_0)\}. \quad (12.10)$$

This approach fully accounts for the basis risk and is perfectly tailored to the risk aversion characteristics of the person or institution facing the risk.

**Example 12.14 (Mean–variance hedging)** One obvious choice for the utility function is the quadratic function

$$U(x) = x - \frac{b}{2}x^2$$

with  $b > 0$ . Then (12.10) leads to a maximization problem involving the means, variances, and covariances of the variables. Smoother derivations and neater formulas are obtained, however, by recognizing that this is essentially equivalent to maximizing the expression

$$V(x) = E(x) - r \operatorname{var}(x)$$

for some positive constant  $r$ . The function  $V$  can be thought of as an altered mean–variance utility.

For meaningful results, the magnitude of  $r$  must be determined by the problem itself. One reasonable choice is  $r = 1/(2\hat{x})$ , where  $\hat{x}$  is a rough estimate of the final value of  $E(x)$ . This then weights variance and one-half of  $[E(x)]^2$  about equally.

Using  $V(x)$  as the objective, the optimal hedging problem becomes

$$\underset{h}{\text{maximize}} \quad \{E[x + h(F_T - F_0)] - r \operatorname{var}(x + hF_T)\}. \quad (12.11)$$

This leads directly, after some algebra, to the solution

$$h = \frac{\bar{F}_T - F_0}{2r \operatorname{var}(F_T)} - \frac{\operatorname{cov}(x, F_T)}{\operatorname{var}(F_T)}. \quad (12.12)$$

Note that the second term is exactly the minimum-variance solution. The first term augments this by accounting for the expected gain due to futures participation. In

---

<sup>6</sup> Ideally, we should express utility in terms of total wealth; but we may assume here that the additional wealth simply changes the definition of  $U$ .

other words, the second term is a pure hedging term, whereas the first term accounts for the fact that hedging is a form of investment, and the expected return of that investment should be incorporated into the portfolio.

This simple formula illustrates, however, the practical difficulty associated with optimal hedging. It is quite difficult to obtain meaningful estimates of  $\bar{F}_T - F_0$ . In fact, in many cases a reasonable estimate is that this difference is zero, so it is understandable why many hedgers prefer to use only the minimum-variance portion of the solution.

**Example 12.15 (The wheat hedge)** Consider the producer of flour and bread of Example 12.8. It is likely that this producer, being a large player in the market, has a good knowledge of wheat market conditions. Suppose that this producer expects the price of wheat to increase by 5% in 3 months. However, the producer recognizes that the wheat market has approximately 30% volatility (per year), so the producer assigns a 15% variation to the 3-month forecast ( $15\% = 30\%/\sqrt{4}$ ).

Using  $x = 500,000F_T$  and applying (12.12), we find

$$\begin{aligned} h &= -500,000 + \frac{1}{2r\text{var}(F_T)}(\bar{F}_T - F_0) \\ &= -500,000 + \frac{1}{2rF_0\text{var}(F_T/F_0)}\left(\frac{\bar{F}_T}{F_0} - 1\right) \\ &= -500,000 + \frac{1}{6.60(.15)^2 r} \times .05 \\ &= -500,000 + \frac{.336}{r}. \end{aligned}$$

Note that the term  $-500,000$  represents the equal and opposite position of perfect hedging. This is augmented by a speculative term, determined by a rough estimate of return on the futures price, and by the value of  $r$ .

Using the method of selecting  $r$  suggested earlier, we have  $r = 1/1,000,000$ . Hence the final hedge is  $h = -500,000 + 336,000 = -164,000$ .

## 12.12 Hedging Nonlinear Risk\*

In our examples so far the risk being hedged was linear, in the sense that final wealth  $x$  was a linear function of an underlying market variable, such as a commodity price. The general theory of hedging does not depend on this assumption, and indeed nonlinear risks frequently occur. For example, immunization of a bond portfolio with T-bills (see Exercise 17) is a nonlinear hedging problem—because the change in the value of a bond portfolio is a nonlinear function of the future T-bill price.

Nonlinear risk can arise in complex contracts. For example, suppose a U.S. firm is negotiating to sell a commodity to a Japanese company at a future date for a price specified in Japanese yen. Both parties recognize that the U.S. firm would face exchange rate risk. Hence an agreement might be made where the U.S. firm

absorbs adverse rate changes up to 10%, while beyond that the two companies share the impact equally.

Nonlinear risks also arise when the price of a good is influenced by the quantity being bought or sold. This situation occurs in farming when the magnitudes of all farmers' crops are mutually correlated, and hence any particular farmer finds that his harvest size is correlated to the market price. We give a detailed example of this type.

**Example 12.16 (A corn farmer)** A certain commodity, which we call corn, is grown by many farmers, but the amount of corn harvested by every farmer depends on the weather: sunny weather yields more corn than cloudy weather during the growing season. All corn is harvested simultaneously, and the price per bushel is determined by a market demand function, which is shown in Figure 12.6. This demand function is

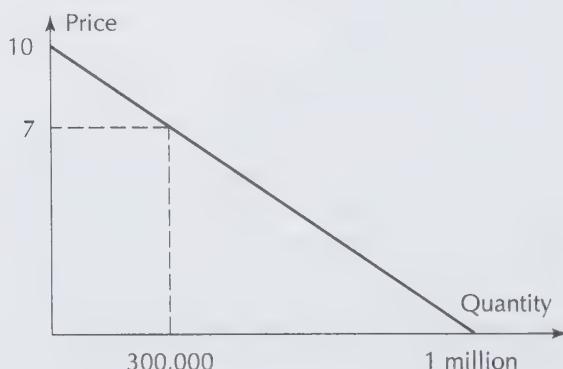
$$P = 10 - D/100,000$$

where  $D$  is the demand (which is also equal, through supply and demand equality, to the total crop size). Each farmer's crop will produce an amount of corn  $C$  which is random. We assume that the amount of corn grown on each farm can vary between 0 and 6,000 bushels, with expected value  $\bar{C} = 3,000$ . The amounts produced on different farms are all perfectly correlated. There are a total of 100 farms, and thus  $D = 100C$  and  $\bar{D} = 300,000$ . The revenue to a farmer will be

$$R = PC = \left(10 - \frac{D}{100,000}\right)C = 10C - \frac{C^2}{1,000}. \quad (12.13)$$

This shows that the revenue is a nonlinear function of the underlying uncertain variable  $C$ . Since  $C$  is random, each farmer faces nonlinear risk.

Can a farmer hedge this risk in advance by participating in the futures market for corn? Try to think this through before we present the analysis. Since the farmer is ultimately going to sell his corn harvest at the (risky) spot price, it might be prudent to



**FIGURE 12.6 Demand for corn.** The price of corn varies from \$10 to \$0 per bushel, depending on the total quantity produced.

**TABLE 12.6**  
**REVENUE FROM PRODUCTION AND HEDGING**

Futures position	Corn production (in 100's of bushels)									
	10	15	20	25	30	35	40	45	50	
50	19000	20250	21000	21250	21000	20250	19000	17250	15000	
45	18000	19500	20500	21000	21000	20500	19500	18000	16000	
40	17000	18750	20000	20750	21000	20750	20000	18750	17000	
35	16000	18000	19500	20500	21000	21000	20500	19500	18000	
30	15000	17250	19000	20250	21000	21250	21000	20250	19000	
25	14000	16500	18500	20000	21000	21500	21500	21000	20000	
20	13000	15750	18000	19750	21000	21750	22000	21750	21000	
15	12000	15000	17500	19500	21000	22000	22500	22500	22000	
10	11000	14250	17000	19250	21000	22250	23000	23250	23000	
5	10000	13500	16500	19000	21000	22500	23500	24000	24000	
0	9000	12750	16000	18750	21000	22750	24000	24750	25000	
-5	8000	12000	15500	18500	21000	23000	24500	25500	26000	
-10	7000	11250	15000	18250	21000	23250	25000	26250	27000	
-15	6000	10500	14500	18000	21000	23500	25500	27000	28000	
-20	5000	9750	14000	17750	21000	23750	26000	27750	29000	
-25	4000	9000	13500	17500	21000	24000	26500	28500	30000	
-30	3000	8250	13000	17250	21000	24250	27000	29250	31000	
-35	2000	7500	12500	17000	21000	24500	27500	30000	32000	
-40	1000	6750	12000	16750	21000	24750	28000	30750	33000	
-45	0	6000	11500	16500	21000	25000	28500	31500	34000	
-50	-1000	5250	11000	16250	21000	25250	29000	32250	35000	

Revenue can be calculated for various futures positions and production outcomes using a spreadsheet.

sell some corn now at a known price in the futures market. Indeed, if the farmer knew exactly how much corn he would produce, and only the price were uncertain, he could implement an equal and opposite policy by shorting this amount in the corn futures market. Perhaps in this actual situation where both amount and price are uncertain, he should short some lesser amount. What do you think?

The way to find the best hedge is to work out the relationships between revenue, production, and the futures position. We assume for simplicity that interest rates are zero. If each farm produces the expected value of  $\bar{C} = 3,000$ , then  $D = 300,000$  and we find  $P = \$7$  per bushel. Hence \$7 represents a nominal anticipated price. Let us assume that \$7 is also the current futures price  $P_0$ . We want to determine the best futures participation.

Let  $h$  be the futures market position. With this position the farmer's revenue will be

$$R = PC + h(P - P_0).$$

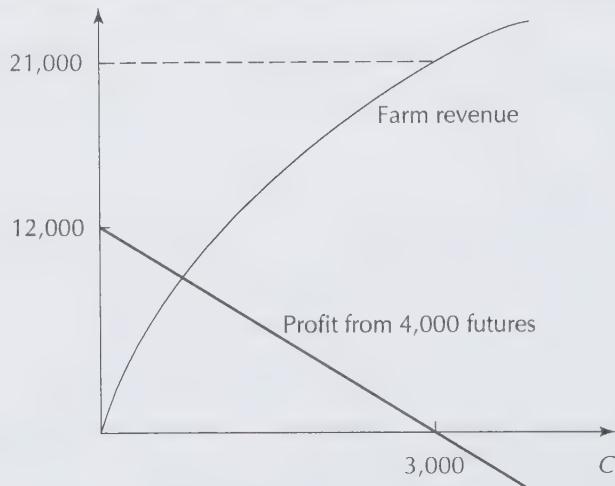
Substituting for  $P$  in terms of  $C$ , we find

$$R = 10C - \frac{C^2}{1,000} + \frac{\bar{C} - C}{1,000}h.$$

This is the equation that the farmer should consider. One simple way to study this equation is to display it in a spreadsheet array, as shown in Table 12.6. This table has the farm's production of corn across the columns, and the futures position (in hundreds of bushels) running along the rows. The entries are the corresponding revenues. For example, note that if the final production is 3,000 bushels (the expected value), then the revenue is \$21,000, independent of the futures position. This is because the final price will be \$7, which is equal to the current futures price; hence the futures contract makes no profit or loss.

The equal and opposite hedge would correspond to a futures position of  $-3,000$  (or  $-30$  in the table). Note that this is actually a very risky position—much more so than the zero position—for the revenue varies widely from \$3,000 to \$31,000. What is the least risky position? We find that position by scanning the rows, looking for the row with the least variation. It is the row marked 40, corresponding to a position of  $+4,000$ . Wow! The optimal position has a sign opposite to that which we might have expected, and a magnitude much greater than the expected value of the crop.<sup>7</sup>

How can we understand the nature of this solution? The original revenue function (12.13), is shown in Figure 12.7. Also shown in the figure is the profit from a  $+4,000$  futures position as a function of the amount of corn grown. Note that the profit from the futures contract *decreases* as more corn is grown. This is because as



**FIGURE 12.7 Farm revenue and hedging.** The best futures position is obtained when the slope of its payoff is equal to and opposite the slope of the revenue.

<sup>7</sup> It can be shown that this position is indeed optimal for any concave increasing utility function if the probabilities of different-size harvests are symmetric. (See Exercise 18.)

more corn is grown, the final spot price of corn decreases. The revenue increases as more corn is produced (although eventually the revenue curve bends downward). At the nominal value of  $C = 3,000$ , the slopes of the two functions are exactly opposite—the slope of the revenue curve is 4 and the slope of the futures profit line is  $-4$ . The two slopes cancel, meaning that the net revenue curve is flat at the nominal point. This is the best linear approximation to the nonlinear hedging problem. Note that the best solution for an individual is not necessarily an equilibrium for all.

Here is one way to think about the situation, to resolve the apparent conundrum. The farmer has a natural hedge against price movements. If the price of corn should go down, the farmer's revenue from corn will go *up* because of his increased harvest, instead of down as it would if the harvest were unaffected. This natural hedge is, in fact, of greater magnitude than an equal and opposite hedge, which would keep net revenue constant. Hence the farmer must counteract the natural hedge by taking a positive position in the futures market.

## 12.13 Summary

A forward contract is a contract to buy or sell an asset at a fixed date in the future. The intrinsic value of a forward contract may vary from day to day, but there are no cash flows until the delivery date. A futures contract is similar, except that it is marked to market daily, with profits or losses flowing to a margin account so that the contract continues to have zero value. The price of a forward contract, in the absence of carrying costs and assuming that the commodity can be shorted, is just  $F = S/d$ , where  $S$  is the current value of the asset and  $d$  is the discount rate that applies for the interval of time until delivery. In other words,  $F$  is the future value of the current spot price  $S$ . If there are carrying costs,  $F$  is the future value of these costs plus the future value of  $S$ . If shorting is not possible, as is frequently the case for consumption assets, the forward price is restricted only to be less than  $S/d$ . For investment assets, the possibility of selling directly can replace the need to sell short.

If interest rates follow expectation dynamics, the prices of a forward contract and a corresponding futures contract are identical, even though their cash flow patterns are slightly different. For analysis purposes, a futures contract can therefore be approximated by the corresponding forward contract.

Forwards and futures are used to hedge risk in commercial transactions. The simplest type of hedge is the perfect, or equal and opposite, hedge, where an obligation to buy or sell a commodity in a future spot market is essentially executed early at a known price by entering a futures contract for the same quantity. If there is no hedging instrument available that matches the commodity of the obligation exactly, a minimum-variance hedge can be constructed using instruments that are correlated with the obligation. A relatively high correlation is required, however, to produce a significant hedging effect.

More sophisticated hedging is obtained by taking an optimal portfolio viewpoint, maximizing expected utility subject to the constraints implied by obligations

and market conditions. This approach has the advantage that it can handle essentially any situation, even those where the decisions affect portfolio value nonlinearly, but it has the disadvantage that detailed information is required. In any case, futures market participation is an important aspect of many hedging operations.

## Exercises

- (Gold futures) The current price of gold is \$412 per ounce. The storage cost is \$2 per ounce per year, payable quarterly in advance. Assuming a constant interest rate of 9% compounded quarterly, what is the theoretical forward price of gold for delivery in 9 months?
- (Proportional carrying charges ◊) Suppose that a forward contract on an asset is written at time zero and there are  $M$  periods until delivery. Suppose that the proportional carrying charge in period  $k$  is  $qS(k)$ , where  $S(k)$  is the spot price of the asset in period  $k$ . Show that the forward price is

$$F = \frac{(1-q)^{-M} S}{d(0, M)}.$$

[Hint: Consider a portfolio that pays all carrying costs by selling a fraction of the asset as required. Let the number of units of the asset held at time  $k$  be  $x(k)$  and find  $x(M)$  in terms of  $x(0)$ .]

- (Silver contract) At the beginning of April one year, the silver forward prices (in cents per troy ounce) were as follows:

Apr	406.50
July	416.64
Sept	423.48
Dec	433.84

(Assume that contracts settle at the end of the given month.) The carrying cost of silver is about 20 cents per ounce per year, paid at the beginning of each month. Estimate the interest rate at that time.

- (Continuous-time carrying charges) Suppose that a continuous-time compounding framework is used with a fixed interest rate  $r$ . Suppose that the carrying charge per unit of time is proportional to the spot price; that is, the charge is  $qS(t)$ . Show that the theoretical forward price of a contract with delivery date  $T$  is

$$F = S e^{(r+q)T}.$$

[Hint: Use Exercise 2.]

- (Carrying cost proof) Complete the second half of the proof of the “forward price formula with carrying cost” in Section 12.3. To construct the arbitrage, go long one unit of a forward and short one unit spot. To execute the short, it is necessary to borrow the asset from someone, say, Mr. X. As part of our arrangement with Mr. X we ask that he give us the carrying costs as they would normally occur, since he would have to pay them if we did not borrow the asset. We then invest these cash flows. At the final time we buy one unit as obligated by our forward and repay Mr. X. Show the details of this argument.

6. (Foreign currency alternative) Consider the situation of Example 12.10. Rather than shorting a futures contract, the U.S. firm could borrow  $500/(1+r_G)$  euros (where  $r_G$  is the 90-day interest rate in Germany), sell these euros into dollars, invest the dollars in T-bills, and then later repay the euros loan with the payment received for the German order. Discuss how this procedure is related to the original one.
7. (A bond forward) A certain 10-year bond is currently selling for \$920. A friend of yours owns a forward contract on this bond that has a delivery date in 1 year and a delivery price of \$940. The bond pays coupons of \$80 every 6 months, with one due 6 months from now and another just before maturity of the forward. The current interest rates for 6 months and 1 year (compounded semiannually) are 7% and 8%, respectively (annual rates compounded every 6 months). What is the current value of the forward contract?
8. (Simple formula) Derive formula (12.6) by converting a cash flow of a bond to that of the fixed portion of the swap.
9. (Integral payoff) Suppose at time 0 you have arranged to be paid at time  $T$  the amount  $\int_0^T S(t)dt$ , where  $S(t)$  is the spot price at  $t$  of a commodity that can be shorted and has zero carrying cost. The interest rate is  $r$ . How much is this arrangement worth now, at time 0?
10. (Currency forward) The interest rates in the UK and the United States are, respectively, 4% and 6% per annum compounded continuously. The spot price of the UK pound is \$1.6. The forward price for a UK pound deliverable in 6 months is \$2.0.
- (a) Determine if an arbitrage opportunity exists.
  - (b) If there is such an opportunity, describe it in detail, showing the risk-free profit.
11. (Equity swap ◊) Mr. A. Gaylord manages a pension fund and believes that his stock selection ability is excellent. However, he is worried because the market could go down. He considers entering an equity swap where each quarter  $i$ , up to quarter  $M$ , he pays counterparty B the previous quarter's total rate of return  $r_i$  on the S&P 500 index times some notional principal and receives payments at a fixed rate  $r$  on the same principal. The total rate of return includes dividends. Specifically,  $1+r_i = (S_i + d_i)/S_{i-1}$ , where  $S_i$  and  $d_i$  are the values of the index at  $i$  and the dividends received from  $i-1$  to  $i$ , respectively. Derive the value of such a swap by the following steps:
- (a) Let  $V_{i-1}(S_i + d_i)$  denote the value at time  $i-1$  of receiving  $S_i + d_i$  at time  $i$ . Argue that  $V_{i-1}(S_i + d_i) = S_{i-1}$  and find  $V_{i-1}(r_i)$ .
  - (b) Find  $V_0(r_i)$ .
  - (c) Find  $\sum_{i=1}^M V_0(r_i)$ .
  - (d) Find the value of the swap.
12. (Forward vanilla) The floating rate portion of a plain vanilla interest rate swap with yearly payments and a notional principal of one unit has cash flows at the end of each year defining a stream starting at time 1 of  $(c_0, c_1, c_2, \dots, c_{M-1})$ , where  $c_i$  is the actual spot rate at the beginning of year  $i$ . Using the concepts of forwards, argue that the value at time zero of  $c_i$  to be received at time  $i+1$  is  $d(0, i+1)r_i$ , where  $r_i$  is the short rate for time  $i$  implied by the current (time zero) term structure and  $d(0, i+1)$  is the implied discount factor to time  $i+1$ . The value of the stream is therefore  $\sum_{i=0}^{M-1} d(0, i+1)r_i$ . Show that this reduces to the formula for  $V$  at the end of Section 12.5.

13. (Specific vanilla) Suppose the current term structure of interest rates is (.070, .073, .077, .081, .084, .088). A plain vanilla interest rate swap will make payments at the end of each year equal to the floating short rate that was posted at the beginning of that year. A 6-year swap having a notional principal of \$10 million is being configured.
- What is the value of the floating rate portion of the swap?
  - What rate of interest for the fixed portion of the swap would make the two sides of the swap equal?
14. (Derivation) Derive the mean-variance hedge formula given by (12.12).
15. (Grapefruit hedge) Farmer D. Jones has a crop of grapefruit that will be ready for harvest and sale as 150,000 pounds of grapefruit juice in 3 months. Jones is worried about possible price changes, so he is considering hedging. There is no futures contract for grapefruit juice, but there is a futures contract for orange juice. His son, Gavin, recently studied minimum-variance hedging and suggests it as a possible approach. Currently the spot prices are \$1.20 per pound for orange juice and \$1.50 per pound for grapefruit juice. The standard deviation of the prices of orange juice and grapefruit juice is about 20% per year, and the correlation coefficient between them is about .7. What is the minimum-variance hedge for farmer Jones, and how effective is this hedge as compared to no hedge?
16. (Opposite hedge variance) Assume that cash flow is given by  $y = S_T W + (F_T - F_0)h$ . Let  $\sigma_s^2 = \text{var}(S_T)$ ,  $\sigma_F^2 = \text{var}(F_T)$ , and  $\sigma_{ST} = \text{cov}(S_T, F_T)$ .
- In an equal and opposite hedge,  $h$  is taken to be an opposite equivalent dollar value of the hedging instrument. Therefore  $h = -kW$ , where  $k$  is the price ratio between the asset and the hedging instrument. Express the standard deviation of  $y$  with the equal and opposite hedge in the form
- $$\sigma_y = W\sigma_S \times B.$$
- (That is, find  $B$ .)
- Apply this to Example 12.12 and compare with the minimum-variance hedge.
17. (Immunization as hedging  $\diamond$ ) A pension fund has just paid some of its liabilities, and as a result of this transaction the fund is no longer fully immunized. The fund manager decides that instead of changing the portfolio, the firm should hedge its position using a futures contract on a Treasury bond. The fund manager wants to hedge against parallel changes to the spot rate curve. Use the following set of information to determine the numerical values of the hedging position:
- Yearly spot rate sequence: .05, .053, .056, .058, .06, .061.
  - Liabilities: \$1 million in 1 year, \$2 million in 2 years, and \$1 million in 3 years.
  - Current bond portfolio: \$4.253 million in par value of zero-coupon bonds maturing in 2 years. (Use the continuous-time formulas for discounting:  $e^{-rt}$ .)
  - The hedge is to be constructed using futures contracts on zero-coupon bonds maturing in 6 years, with a contract delivery date in 1 year.
18. (Symmetric probability  $\diamond$ ) Suppose the wealth that is to be received at a time  $T$  in the future has the form

$$W = a + hx + cx^2,$$

where  $a$  is a constant and  $x$  is a random variable. The value of the variable  $h$  can be selected by the investor. Suppose that the investor has a utility function that is increasing and strictly concave. Suppose also that the probability distribution of  $x$  is symmetric; that is,  $x$  and  $-x$

have the same distribution. It follows that  $E(x) = 0$  and that the investor cannot influence the expected value of wealth.

- (a) Show that the optimal choice is  $h = 0$ .
- (b) Apply this result to the corn farm problem to show that the optimal futures position is +4,000.

**19.** (Double symmetric probability ◊) Suppose that revenue has the form

$$R = Axy + Bx - hy,$$

where  $h$  can be chosen and  $x$  and  $y$  are random variables. The distribution of  $x$  and  $y$  is symmetric about  $(0,0)$ ; that is,  $-x, -y$  has the same distribution as  $x, y$ . Show that the choice of  $h$  that minimizes the variance of  $R$  is

$$h = B\sigma_{xy}/\sigma_y^2.$$

**20.** (A general farm problem ◊) Suppose that, as in the corn farm example, the farm has random production and the final spot price is governed by the same demand function. However, the crop of the farm is not perfectly correlated to total demand, but  $\sigma_{CD}$  and  $\sigma_D^2$  are known. The current futures price is also equal to the expected final spot price. Show that the minimum-variance hedging position is

$$h = 100,000 \left( \frac{-3}{100} + \frac{7\sigma_{CD}}{\sigma_D^2} \right).$$

Check the solution for the special cases (a)  $D = 100C$  and (b)  $\sigma_{CD} = 0$ . [Hint: Use Exercise 19.]

## References

There are several books devoted to futures markets; for example, [1–3]. An excellent book, similar in level to this textbook, is [4]. The futures-forward equivalence result was proved in [5] for the case of a constant interest rate. See [6] for a discussion of hedging techniques, and [7] for the use of interest rate futures similar to that of Exercise 17.

1. Duffie, D. (1989), *Futures Markets*, Prentice Hall, Englewood Cliffs, NJ.
2. Teweles, R. J., and F. J. Jones (1987), *The Futures Game*, McGraw-Hill, New York.
3. Stoll, H. R., and R. E. Whaley (1993), *Futures and Options*, South-West Publishing, Cincinnati, OH.
4. Hull, J. C. (2008), *Options, Futures, and Other Derivative Securities*, 7th ed., Prentice Hall, Englewood Cliffs, NJ.
5. Cox, J. C., J. E. Ingersoll, and S. A. Ross (1981), “The Relation between Forward Prices and Futures Prices,” *Journal of Financial Economics*, 9, 321–346.
6. Figlewski, S. (1986), *Hedging with Financial Futures for Institutional Investors*, Ballenger Publishing, Cambridge, MA.
7. Kolb, R. W., and G. D. Gay (1982), “Immunizing Bond Portfolios with Interest Rate Futures,” *Financial Management*, 11, 81–89.

## MODELS OF ASSET DYNAMICS

True multiperiod investments fluctuate in value, distribute random dividends, exist in an environment of variable interest rates, and are subject to a continuing variety of other uncertainties. This chapter initiates the study of such investments by showing how to model asset price fluctuations conveniently and realistically. This chapter therefore contains no investment principles as such. Rather it introduces the mathematical models that form the foundation for the analyses developed in later chapters.

Two primary model types are used to represent asset dynamics: binomial lattices and Ito processes. Binomial lattices are analytically simpler than Ito processes, and they provide an excellent basis for computational work associated with investment problems. For these reasons it is best to study binomial lattice models first. The important investment concepts can all be expressed in terms of these models, and many real investment problems can be formulated and solved using the binomial lattice framework. Indeed, roughly 80% of the material in later chapters is presented in terms of binomial lattice models.

Ito processes are more realistic than binomial lattice models in the sense that they have a continuum of possible stock prices at each period, not just two. Ito process models also allow some problems to be solved analytically, as well as computationally. They also provide the foundation for constructing binomial lattice models in a clear and consistent manner. For these reasons Ito process models are fundamental to dynamic problems. For a complete understanding of investment principles, it is important to understand them.

The organization of this chapter is based on the preceding viewpoint concerning the roles of different models. The first section presents the binomial lattice model directly. With this background, much of the material in later chapters can be studied.

The remaining sections consider models that have a continuum of price values. These models are developed progressively from discrete-time models to continuous-time models based on Ito processes.

## 13.1 Binomial Lattice Model

To define a binomial lattice model, a basic period length is established (such as 1 week). According to the model, if the price is known at the beginning of a period, the price at the beginning of the next period is one of only two possible values. Usually these two possibilities are defined to be multiples of the price at the previous period—a multiple  $u$  (for up) and a multiple  $d$  (for down). Both  $u$  and  $d$  are positive, with  $u > 1$  and (usually)  $d < 1$ . Hence if the price at the beginning of a period is  $S$ , it will be either  $uS$  or  $dS$  at the next period. The probabilities of these possibilities are  $p$  and  $1 - p$ , respectively, for some given probability  $p$ ,  $0 < p < 1$ . That is, if the current price is  $S$ , there is a probability  $p$  that the new price will be  $uS$  and a probability  $1 - p$  that it will be  $dS$ . This model continues on for several periods.

The general form of such a lattice is shown in Figure 13.1. The stock price can be visualized as moving from node to node in a rightward direction. The probability of an upward movement from any node is  $p$  and the probability of a downward movement is  $1 - p$ . A lattice is the appropriate structure in this case, rather than a tree, because an up movement followed by a down is identical to a down followed by an up. Both produce  $ud$  times the price.

The model may at first seem too simple because it permits only two possible values at the next period. But if the period length is small, many values are possible after several short steps.

To specify the model completely, we must select values for  $u$  and  $d$  and the probability  $p$ . These should be chosen in such a way that the true stochastic nature of the stock is captured as faithfully as possible, as will be discussed.

Because the model is multiplicative in nature (the new value being  $uS$  or  $dS$ , with  $u > 0, d > 0$ ), the price will never become negative. It is therefore possible to consider the logarithm of price as a fundamental variable. For reasons discussed in later sections, use of the logarithm is in fact very helpful and leads to simple formulas for selecting the parameters.

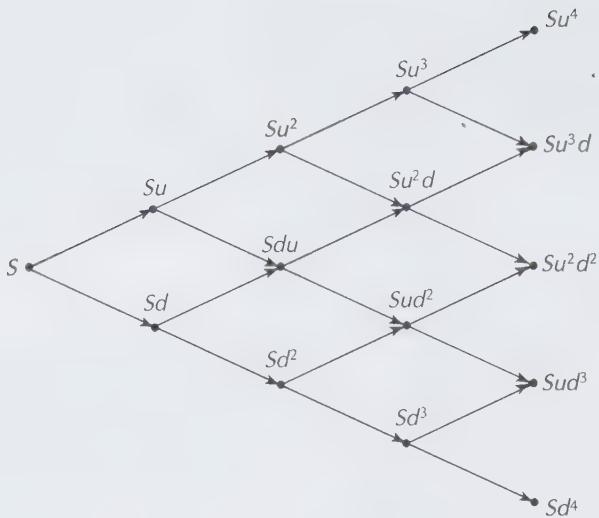
Accordingly, we define  $v$  as the expected yearly growth rate.<sup>1</sup> Specifically,

$$v = E[\ln(S_T/S_0)],$$

where  $S_0$  is the initial stock price and  $S_T$  is the price at the end of 1 year.

---

<sup>1</sup> If the process were deterministic, then  $v = \ln(S_T/S_0)$  implies  $S_T = S_0 e^{vT}$ , which shows that  $v$  is the exponential growth rate.



**FIGURE 13.1 Binomial lattice stock model.** At each step the stock price  $S$  either increases to  $uS$  or decreases to  $dS$ .

Likewise, we define  $\sigma$  as the yearly standard deviation. Specifically,

$$\sigma^2 = \text{var}[\ln(S_T/S_0)].$$

If a period length of  $\Delta t$  is chosen, which is small compared to 1, the parameters of the binomial lattice can be selected as

$$\begin{aligned} p &= \frac{1}{2} + \frac{1}{2} \left( \frac{\nu}{\sigma} \right) \sqrt{\Delta t} \\ u &= e^{\sigma \sqrt{\Delta t}} \\ d &= e^{-\sigma \sqrt{\Delta t}}. \end{aligned} \tag{13.1}$$

With this choice, the binomial model will closely match the values of  $\nu$  and  $\sigma$  (as shown later); that is, the expected growth rate of  $\ln S$  in the binomial model will be nearly  $\nu$ , and the variance of that rate will be nearly  $\sigma^2$ . The closeness of the match improves if  $\Delta t$  is made smaller, becoming exact as  $\Delta t$  goes to zero.

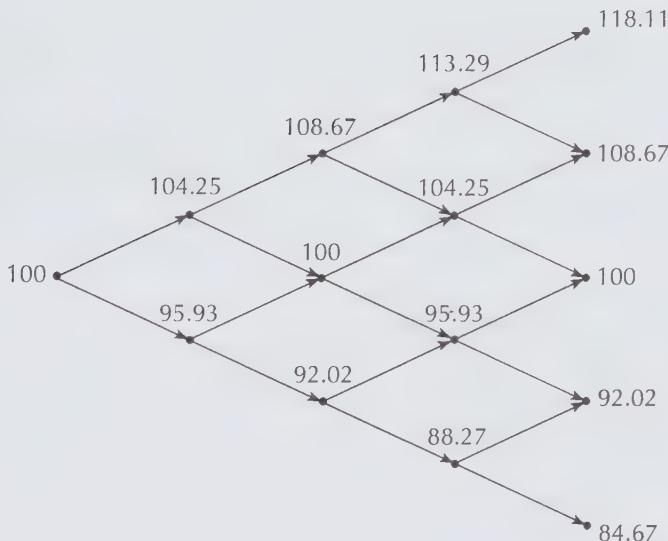
**Example 13.1 (A volatile stock)** Consider a stock with the parameters  $\nu = 15\%$  and  $\sigma = 30\%$ . We wish to make a binomial model based on weekly periods. According to (13.1), we set

$$u = e^{30/\sqrt{52}} = 1.04248, \quad d = 1/u = .95925,$$

and

$$p = \frac{1}{2} \left( 1 + \frac{15}{30} \sqrt{\frac{1}{52}} \right) = .534669.$$

The lattice for this example is shown in Figure 13.2, assuming  $S(0) = 100$ .



**FIGURE 13.2 Lattice for Example 13.1.** The parameters are chosen so that the expected growth rate of the logarithm of price and the variance of that growth rate match the known corresponding values for the asset.

We shall return to the binomial lattice later in this chapter after studying models that allow a continuum of prices. The binomial model will be found to be a natural approximation to these models.

## 13.2 The Additive Model

We now study models with the property that price can range over a continuum. First we shall consider discrete-time models, beginning with the additive model of this section, and then later we shall consider continuous-time models defined by Ito processes.

Let us focus on  $N + 1$  time points, indexed by  $k$ ,  $k = 0, 1, 2, \dots, N$ . We also focus on a particular asset that is characterized by a price at each time. The price at time  $k$  is denoted by  $S(k)$ . Our model will recognize that the price in any one time is dependent to some extent on previous prices.

The simplest model is the **additive model**,

$$S(k+1) = aS(k) + u(k) \quad (13.2)$$

for  $k = 0, 1, 2, \dots, N - 1$ . In this equation,  $a$  is a constant (usually  $a > 1$ ) and the quantities  $u(k)$ ,  $k = 0, 1, \dots, N - 1$ , are random variables. The  $u(k)$ 's can be thought of as “shocks” or “disturbances” that cause the price to fluctuate. To operate or run

this model, an initial price  $S(0)$  is specified; then once the random variable  $u(0)$  is given,  $S(1)$  can be determined. The process then repeats progressively in a stepwise fashion, determining  $S(2), S(3), \dots, S(N)$ .

The key ingredient of this model is the sequence of random variables  $u(k)$ ,  $k = 1, 2, \dots, N$ . We assume that these are mutually statistically independent.

Note that the price at any time depends only on the price at the most recent previous time and the random disturbance. It does not explicitly depend on other previous prices.

## Normal Price Distribution

It is instructive to solve explicitly for a few of the prices from (13.2). By direct substitution we have

$$\begin{aligned} S(1) &= aS(0) + u(0) \\ S(2) &= aS(1) + u(1) \\ &= a^2S(0) + au(0) + u(1). \end{aligned}$$

By simple induction it can be seen that for general  $k$ ,

$$S(k) = a^k S(0) + a^{k-1}u(0) + a^{k-2}u(1) + \dots + u(k-1). \quad (13.3)$$

Hence  $S(k)$  is  $a^k S(0)$  plus the sum of  $k$  random variables.

Frequently we assume that the random variables  $u(k)$ ,  $k = 0, 1, 2, \dots, N - 1$ , are independent normal random variables with a common variance  $\sigma^2$ . Then, since a linear combination of normal random variables is also normal (see Appendix A), it follows from (13.3) that  $S(k)$  is itself a normal random variable.

If the expected values of all the  $u(k)$ 's are zero, then the expected value of  $S(k)$  is

$$E[S(k)] = a^k S(0).$$

When  $a > 1$ , this model has the property that the expected value of the price increases geometrically (that is, according to  $a^k$ ). Indeed, the constant  $a$  is the growth rate factor of the model.

The additive model is structurally simple and easy to work with. The expected value of price grows geometrically, and all prices are normal random variables. However, the model is seriously flawed because it lacks realism. Normal random variables can take on negative values, which means that the prices in this model might be negative as well; but real stock prices are never negative. Furthermore, if a stock were to begin at a price of, say, \$1 with a  $\sigma$  of, say, \$.50 and then drift upward to a price of \$100, it seems very unlikely that the  $\sigma$  would remain at \$.50. It is more likely that the standard deviation would be proportional to the price. For these reasons the additive model is not a good general model of asset dynamics. The model is useful for localized analyses, over short periods of time (perhaps up to a few months for common stocks), and it is a useful building block for other models, but it cannot be used alone as an ongoing model representing long- or intermediate-term fluctuations. For

this reason we must consider a better alternative, which is the multiplicative model. (However, our understanding of the additive model will be important for that more advanced model.)

## 13.3 The Multiplicative Model

The **multiplicative model** has the form

$$S(k+1) = u(k)S(k) \quad (13.4)$$

for  $k = 0, 1, 2, \dots, N - 1$ . Here again the quantities  $u(k)$ ,  $k = 0, 1, 2, \dots, N - 1$ , are mutually independent random variables. The variable  $u(k)$  defines the *relative* change in price between times  $k$  and  $k + 1$ . This relative change is  $S(k+1)/S(k)$ , which is independent of the overall magnitude of  $S(k)$ . It is also independent of the units of price. For example, if we change units from U.S. dollars to German marks, the relative price change is still  $u(k)$ .

The multiplicative model takes a familiar form if we take the natural logarithm of both sides of the equation. This yields

$$\ln S(k+1) = \ln S(k) + \ln u(k) \quad (13.5)$$

for  $k = 0, 1, 2, \dots, N - 1$ . Hence in this form the model is of the additive type with respect to the logarithm of the price, rather than the price itself. Therefore we can use our knowledge of the additive model to analyze the multiplicative model.

It is now natural to specify the random disturbances directly in terms of the  $\ln u(k)$ 's. In particular we let

$$w(k) = \ln u(k)$$

for  $k = 0, 1, 2, \dots, N - 1$ , and we specify that these  $w(k)$ 's be normal random variables. We assume that they are mutually independent and that each has expected value  $\bar{w}(k) = v$  and variance  $\sigma^2$ .

We can express the original multiplicative disturbances as

$$u(k) = e^{w(k)} \quad (13.6)$$

for  $k = 0, 1, 2, \dots, N - 1$ . Each of the variables  $u(k)$  is said to be a **lognormal** random variable since its logarithm is in fact a normal random variable.

Notice that now there is no problem with negative values. Although the normal variable  $w(k)$  may be negative, the corresponding  $u(k)$  given by (13.6) is always positive. Since the random factor by which a price is multiplied is  $u(k)$ , it follows that prices remain positive in this model.

### Lognormal Prices

The successive prices of the multiplicative model can be easily found to be

$$S(k) = u(k-1)u(k-2)\cdots u(0)S(0).$$

Taking the natural logarithm of this equation we find

$$\ln S(k) = \ln S(0) + \sum_{i=0}^{k-1} \ln u(i) = \ln S(0) + \sum_{i=0}^{k-1} w(i).$$

The term  $\ln S(0)$  is a constant, and the  $w(i)$ 's are each normal random variables. Since the sum of normal random variables is itself a normal random variable (see Appendix A), it follows that  $\ln S(k)$  is normal. In other words, all prices are lognormal under the multiplicative model.

If each  $w(i)$  has expected value  $\bar{w}(i) = v$  and variance  $\sigma^2$ , and all are mutually independent, then we find

$$E[\ln S(k)] = \ln S(0) + vk \quad (13.7a)$$

$$\text{var}[\ln S(k)] = k\sigma^2. \quad (13.7b)$$

Hence both the expected value and the variance increase linearly with  $k$ .

## Real Stock Distributions

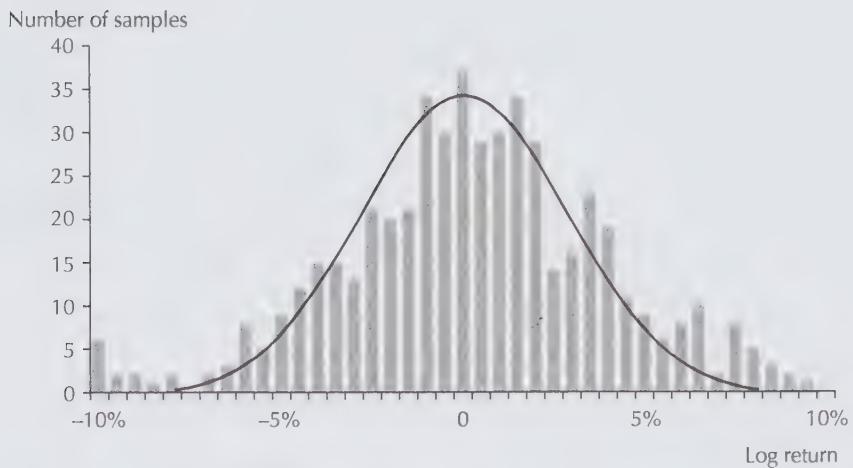
At this point it is natural to ask how well this theoretical model fits actual stock price behavior. Are real stock prices lognormal?

The answer is that, based on an analysis of past stock price records, the price distributions of most stocks are actually quite close to lognormal. To verify this, we select a nominal period length of, say, 1 week and record the differences  $\ln S(k+1) - \ln S(k)$  for many values of  $k$ ; that is, we record the weekly changes in the logarithm of the prices for many weeks. We then construct a histogram of these values and compare it with that of a normal of the same mean and variance.

Typically, the measured histogram is similar to that of a corresponding lognormal density, but it does exhibit some significant differences. These differences are noticeable in the extreme ends—the **tails** of the density. Often the tails account for a greater share of probability than does a lognormal density. This is expressed by saying that the tails are heavier or that the density has “fat tails.” Another common feature of the observed histogram is that one tail may be heavier than the other. The density is then said to be *skewed*. Negative skewness reflects that the lower tail is heavier than the upper.<sup>2</sup> The fact that some stocks exhibit negative skewness is logical when we notice that stock prices sometimes drop quickly but generally recover slowly. Many assets exhibit both heavy tails and negative skewness. (See Figure 13.3.) For many applications, such as design of a portfolio of stocks, these differences are not serious, since the variation in tails represents a small fraction of the entire density. However, some investments, such as “out-of-the-money” stock options, are very much influenced by extreme movements of the underlying stocks.

---

<sup>2</sup> The skewness of log return is  $E[\ln S - E(\ln S)]^3 / \sigma^3$ .



**FIGURE 13.3 Observed density.** Histogram of Walt Disney Inc. weekly log stock returns during 2002–2012. The distribution exhibits fat tails and essentially zero skewness.

## 13.4 Typical Parameter Values\*

The return of a stock over the period between  $k$  and  $k + 1$  is  $S(k+1)/S(k)$ , which under the multiplicative model is equal to  $u(k)$ . The value of  $w(k) = \ln u(k)$  is therefore the logarithm of the return. The mean value of  $w(k)$  is denoted by  $\nu$  and the variance of  $w(k)$  by  $\sigma^2$ . Typical values of these parameters for assets such as common stocks can be inferred from our knowledge of corresponding values for returns. Thus for stocks, typical values of  $\nu = E[w(k)]$  and  $\sigma = \text{stdev}[w(k)]$  might be

$$\nu = 12\%, \quad \sigma = 15\%$$

when the length of a period is 1 year. If the period length is less than a year, these values scale downward;<sup>3</sup> that is, if the period length is  $p$  part of a year, then

$$\nu_p = p\nu, \quad \sigma_p = \sqrt{p}\sigma.$$

The values can be estimated from historical records in the standard fashion (but with caution as to the validity of these estimates, as raised in Chapter 9). If we have  $N + 1$  time points of data, spanning  $N$  periods, the estimate of the single-period  $\nu$  is

$$\begin{aligned} \hat{\nu} &= \frac{1}{N} \sum_{k=0}^{N-1} \ln \left[ \frac{S(k+1)}{S(k)} \right] = \frac{1}{N} \sum_{k=0}^{N-1} [\ln S(k+1) - \ln S(k)] \\ &= \frac{1}{N} \ln \left[ \frac{S(N)}{S(0)} \right]. \end{aligned}$$

---

<sup>3</sup> Using log returns, the scaling is *exactly* proportional. There is no error due to compounding as with returns (without the log). (See Exercise 2.)

Hence all that matters is the ratio of the last to the first price.

The standard estimate of  $\sigma^2$  is

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{k=0}^{N-1} \left\{ \ln \left[ \frac{S(k+1)}{S(k)} \right] - \hat{\nu} \right\}^2.$$

As with the estimation of return parameters, the error in these estimates can be characterized by their variances. For  $\nu$  this variance is

$$\text{var}(\hat{\nu}) = \sigma^2/N,$$

and for  $\sigma^2$  it is [assuming  $w(k)$  is normal]

$$\text{var}(\hat{\sigma}^2) = 2\sigma^4/(N-1).$$

Hence for the values assumed earlier, namely,  $\nu = .12$  and  $\sigma = .15$ , we find that 10 years of data is required to reduce the standard deviation of the estimate<sup>4</sup> of  $\nu$  to .05 (which is still a sizable fraction of the true value). On the other hand, with only 1 year of weekly data we can obtain a fairly good estimate<sup>5</sup> of  $\sigma^2$ .

## 13.5 Lognormal Random Variables

If  $u$  is a lognormal random variable, then the variable  $w = \ln u$  is normal. In this case we found that the prices in the multiplicative model are all lognormal random variables. It is therefore useful to study a few important properties of such random variables.

The general shape of the probability density of a lognormal random variable is shown in Figure 13.4. Note that the variable is always nonnegative and the density is somewhat skewed.

Suppose that  $w$  is normal and has expected value  $\bar{w}$  and variance  $\sigma^2$ . What is the expected value of  $u = e^w$ ? A quick guess might be  $\bar{u} = e^{\bar{w}}$ , but this is wrong. Actually  $\bar{u}$  is greater than this by the factor  $e^{\frac{1}{2}\sigma^2}$ ; that is,

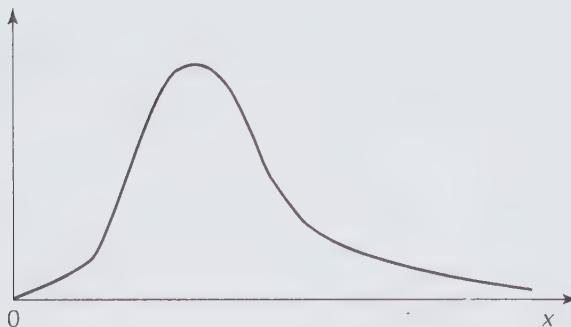
$$\bar{u} = e^{\bar{w} + \frac{1}{2}\sigma^2}. \quad (13.8)$$

This result can be intuitively understood by noting that as  $\sigma$  is increased, the lognormal density will spread out. It cannot spread downward below zero, but it can spread upward unboundedly. Hence the mean value increases as  $\sigma$  increases.

The extra term  $\frac{1}{2}\sigma^2$  is actually fairly small for low-volatility stocks. For example, consider a stock with a yearly  $\bar{w} = .12$  and a yearly  $\sigma$  of .15. The correction term is  $\frac{1}{2}\sigma^2 = .0225$ , which is small compared to  $\bar{w}$ . For stocks with high volatility, however, the correction can be significant.

<sup>4</sup>  $\sigma(\hat{\nu}) = \frac{\sigma}{\sqrt{N}} = \frac{\sigma}{\sqrt{10}} = \frac{.15}{3.16} = .05$ .

<sup>5</sup> See Section 9.2.



**FIGURE 13.4 Lognormal density.** The lognormal probability density is nonzero only for  $x > 0$ .

## 13.6 Random Walks and Wiener Processes

In Section 13.7 we will shorten the period length in a multiplicative model and take the limit as this length goes to zero. This will produce a model in continuous time. In preparation for that step, we introduce special random functions of time, called random walks and Wiener processes.

Suppose that we have  $N$  periods of length  $\Delta t$ . We define the additive process  $z$  by

$$z(t_{k+1}) = z(t_k) + \epsilon(t_k)\sqrt{\Delta t} \quad (13.9)$$

$$t_{k+1} = t_k + \Delta t \quad (13.10)$$

for  $k = 0, 1, 2, \dots, N$ . This process is termed a **random walk**. In these equations  $\epsilon(t_k)$  is a normal random variable with mean 0 and variance 1—a **standardized normal random variable**. These random variables are mutually uncorrelated; that is,  $E[\epsilon(t_j)\epsilon(t_k)] = 0$  for  $j \neq k$ . The process is started by setting  $z(t_0) = 0$ . Thereafter a particular realized path wanders around according to the happenstance of the random variables  $\epsilon(t_k)$ . [The reason for using  $\sqrt{\Delta t}$  in (13.9) will become clear shortly.] A particular path of a random walk is shown in Figure 13.5.

Of special interest are the difference random variables  $z(t_k) - z(t_j)$  for  $j < k$ . We can write such a difference as

$$z(t_k) - z(t_j) = \sum_{i=j}^{k-1} \epsilon(t_i)\sqrt{\Delta t}.$$

This is a normal random variable because it is the sum of normal random variables. We find immediately that

$$E[z(t_k) - z(t_j)] = 0.$$

**FIGURE 13.5 Possible random walk.**

The movements are determined by normal random variables.



Also, using the independence of the  $\epsilon(t_k)$ 's, we find

$$\begin{aligned}\text{var}[z(t_k) - z(t_j)] &= \text{E} \left[ \sum_{i=j}^{k-1} \epsilon(t_i) \sqrt{\Delta t} \right]^2 \\ &= \text{E} \left[ \sum_{i=j}^{k-1} \epsilon(t_i)^2 \Delta t \right] \\ &= (k-j)\Delta t = t_k - t_j.\end{aligned}$$

Hence the variance of  $z(t_k) - z(t_j)$  is exactly equal to the time difference  $t_k - t_j$  between the points. This calculation also shows why  $\sqrt{\Delta t}$  was used in the definition of the random walk so that  $\Delta t$  would appear in the variance.

It should be clear that the difference variables associated with two different time intervals are uncorrelated if the two intervals are nonoverlapping. That is, if  $t_{k_1} < t_{k_2} \leq t_{k_3} < t_{k_4}$ , then  $z(t_{k_2}) - z(t_{k_1})$  is uncorrelated with  $z(t_{k_4}) - z(t_{k_3})$  because each of these differences is made up of different  $\epsilon$ 's, which are themselves uncorrelated.

A **Wiener process** is obtained by taking the limit of the random walk process (13.9) as  $\Delta t \rightarrow 0$ . In symbolic form we write the equations governing a Wiener process as

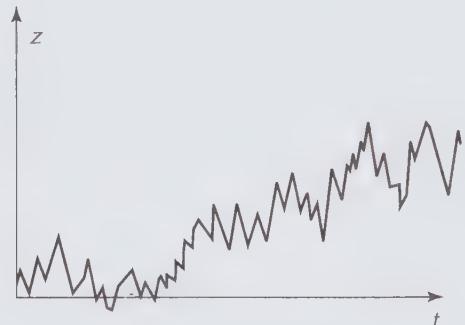
$$dz = \epsilon(t) \sqrt{dt} \quad (13.11)$$

where each  $\epsilon(t)$  is a standardized normal random variable. The random variables  $\epsilon(t')$  and  $\epsilon(t'')$  are uncorrelated whenever  $t' \neq t''$ .

This description of a Wiener process is not rigorous because we have no assurance that the limiting operations are defined; but it provides a good intuitive description. An alternative definition of a Wiener process can be made by simply listing the required properties. In this approach we say a process  $z(t)$  is a **Wiener process** (or, alternatively, **Brownian motion**) if it satisfies the following:

1. For any  $s < t$  the quantity  $z(t) - z(s)$  is a normal random variable with mean zero and variance  $t - s$ .

**FIGURE 13.6 Path of a Wiener process.** A Wiener process moves continuously but is not differentiable.



2. For any  $0 \leq t_1 < t_2 \leq t_3 < t_4$ , the random variables  $z(t_2) - z(t_1)$  and  $z(t_4) - z(t_3)$  are uncorrelated.
3.  $z(t_0) = 0$  with probability 1.

These properties parallel the properties of the random walk process given earlier.

It is fun to try to visualize the outcome of a Wiener process. A sketch of a possible path is shown in Figure 13.6. Remember that given  $z(t)$  at time  $t$ , the value of  $z(s)$  at time  $s > t$  is, on average, the same as  $z(t)$  but will vary from that according to a standard deviation equal to  $\sqrt{s-t}$ .

A Wiener process is not differentiable with respect to time. We can verify this roughly by noting that, for  $t < s$ ,

$$E \left[ \frac{z(s) - z(t)}{s - t} \right]^2 = \frac{s - t}{(s - t)^2} = \frac{1}{s - t} \rightarrow \infty$$

as  $s \rightarrow t$ .

It is, however, useful to have a word for the term  $dz/dt$  since this expression appears in many stochastic equations. A common word used, arising from the systems engineering field (the field that motivated Wiener's work), is **white noise**. It is really fun to try to visualize white noise. One depiction is presented in Figure 13.7.

## Generalized Wiener Processes and Ito Processes

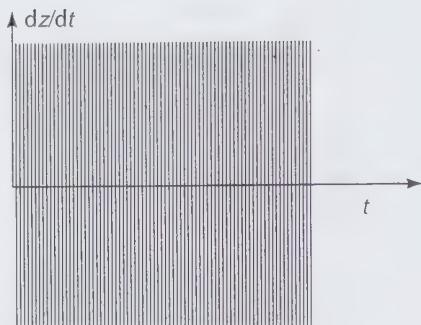
The Wiener process (or Brownian motion) is the fundamental building block for a whole collection of more general processes. These generalizations are obtained by inserting white noise in an ordinary differential equation.

The simplest extension of this kind is the **generalized Wiener process**, which is of the form

$$dx(t) = adt + bdz, \quad (13.12)$$

where  $x(t)$  is a random variable for each  $t$ ,  $z$  is a Wiener process, and  $a$  and  $b$  are constants.

**FIGURE 13.7 Fantasizing white noise.**  
White noise is the derivative of a Wiener process, but that derivative does not exist in the normal sense.



A generalized Wiener process is especially important because it has an analytic solution (which can be found by integrating both sides). Specifically,

$$x(t) = x(0) + at + bz(t). \quad (13.13)$$

An **Ito process** is somewhat more general still. Such a process is described by an equation of the form

$$dx(t) = a(x, t) dt + b(x, t) dz. \quad (13.14)$$

As before,  $z$  denotes a Wiener process. Now, however, the coefficients  $a(x, t)$  and  $b(x, t)$  may depend on  $x$  and  $t$ , and a general solution cannot be written in an analytic form. A special form of Ito process is used frequently to describe the behavior of financial assets, as discussed in the next section.

## 13.7 A Stock Price Process

We now have the tools necessary to extend the multiplicative model of stock prices to a continuous-time model. Recall that the multiplicative model is

$$\ln S(k+1) - \ln S(k) = w(k),$$

where the  $w(k)$ 's are uncorrelated normal random variables. The continuous-time version of this equation is

$$d \ln S(t) = v dt + \sigma dz, \quad (13.15)$$

where  $v$  and  $\sigma \geq 0$  are constants and  $z$  is a standard Wiener process. The whole right-hand side of the equation can be regarded as playing the role of the random variable  $w(k)$  in the discrete-time model. This side can be thought of as a constant plus a normal random variable with zero mean, and hence, overall it is a normal random variable. (Although all terms in the equation are differentials or multiples of differentials and thus do not themselves have magnitude in the usual sense, it is helpful to think of  $dt$  and  $dz$  as being “small” like  $\Delta t$  and  $\Delta z$ .) The term  $vdt$  is, accordingly, the mean value of the right-hand side. This mean value is proportional

to  $dt$ , consistent with the fact that in the logarithm version of the multiplicative model the mean value of the change in  $\ln S$  is proportional to the length of one period. The standard deviation of the right-hand side is  $\sigma$  times the standard deviation of  $dz$ . Hence it is of order of magnitude  $\sigma\sqrt{dt}$ , which is consistent with the fact that in the logarithm version of the multiplicative model the standard deviation of the change in  $\ln S$  is proportional to the square root of the length of one period, as reflected by (13.7a) and (13.7b).

Since equation (13.15) is expressed in terms of  $\ln S(t)$ , it is actually a generalized Wiener process. Hence we can solve it explicitly using (13.13) as

$$\ln S(t) = \ln S(0) + vt + \sigma \dot{z}(t). \quad (13.16)$$

This shows that  $E[\ln S(t)] = E[\ln S(0)] + vt$ , and hence  $E[\ln S(t)]$  grows linearly with  $t$ . Because the expected logarithm of this process increases linearly with  $t$ , just as a continuously compounded bank account, this process is termed **geometric Brownian motion**.

## Lognormal Prices

Like the discrete-time multiplicative model, the geometric Brownian motion process described by (13.15) is a lognormal process. This can be seen easily from the solution (13.16). The right-hand side of that equation is a normal random variable with expected value  $\ln S(0) + vt$  and standard deviation  $\sigma\sqrt{t}$ .

We conclude that the price  $S(t)$  itself has a lognormal distribution. We can express this formally by  $\ln S(t) \sim N(\ln S(0) + vt, \sigma^2 t)$ , where  $N(m, \sigma^2)$  denotes the normal distribution with mean  $m$  and variance  $\sigma^2$ .

Although we can write  $S(t) = \exp[\ln S(t)] = S(0) \exp[v t + \sigma z(t)]$ , it does *not* follow that the expected value of  $S(t)$  is  $S(0)e^{vt}$ . The mean value must instead be determined by equation (13.8), the general formula that applies to lognormal variables. Hence,

$$E[S(t)] = S(0)e^{(v+\frac{1}{2}\sigma^2)t}.$$

If we define  $\mu = v + \frac{1}{2}\sigma^2$ , we have

$$E[S(t)] = S(0)e^{\mu t}.$$

The standard deviation of  $S(t)$  is also given by a general relation for lognormal variables. In the case of the standard deviation, the required calculation is a bit more complex. The formula is (see Exercise 5)

$$\text{stdev}[S(t)] = S(0)e^{vt+\frac{1}{2}\sigma^2 t}(e^{\sigma^2 t} - 1)^{1/2}.$$

## Standard Ito Form

We have defined the random process for  $S(t)$  in terms of  $\ln S(t)$  rather than directly in terms of  $S(t)$ . The use of  $\ln S(t)$  facilitated the development, and it highlights the

fact that the process is a straightforward generalization of the multiplicative model that leads to lognormal distributions. It is, however, useful to express the process in terms of  $S(t)$  itself.

In ordinary calculus we know that

$$d \ln[S(t)] = \frac{dS(t)}{S(t)}.$$

Hence we might be tempted to substitute  $dS(t)/S(t)$  for  $d \ln S(t)$  in the basic equation [Eq. (13.15)], obtaining  $dS(t)/S(t) = vdt + \sigma dz$ . This would be almost correct, but there is a correction term that must be applied when changing variables in Ito processes (because Wiener processes are not ordinary functions and do not follow the rules of ordinary calculus). The appropriate Ito process in terms of  $S(t)$  is

$$\frac{dS(t)}{S(t)} = \left( v + \frac{1}{2}\sigma^2 \right) dt + \sigma dz. \quad (13.17)$$

Note that the correction term  $\frac{1}{2}\sigma^2$  is exactly the same as needed in the expression for the expected value of a lognormal random variable. Putting  $\mu = v + \frac{1}{2}\sigma^2$ , we may write the equation in the standard Ito form for price dynamics,

$$\frac{dS(t)}{S(t)} = \mu dt + \sigma dz. \quad (13.18)$$

The term  $dS(t)/S(t)$  can be thought of as the differential return of the stock; hence in this form the differential return has a simple form.

The correction term required when transforming the equation from  $\ln S(t)$  to  $S(t)$  is a special instance of a general transformation equation defined by **Ito's lemma**, which applies to variables defined by Ito processes. Ito's lemma is discussed in the next section.

Note that if the equation in standard form is written with  $S$  in the denominator, as in (13.17), it is an equation for  $dS/S$ . This term can be interpreted as the instantaneous rate of return on the stock. Hence the standard form is often referred to as an equation for the instantaneous return.

**Example 13.2 (Bond price dynamics)** Let  $P(t)$  denote the price of a bond that pays \$1 at time  $t = T$ , with no other payments. Assume that interest rates are constant at  $r$ . The price of this bond satisfies

$$\frac{dP(t)}{P(t)} = rdt,$$

which is a deterministic Ito equation, paralleling the equation for stock prices. The solution to this equation is  $P(t) = P(0)e^{rt}$ . Using  $P(T) = 1$ , we find that  $P(t) = e^{r(t-T)}$ .

We now summarize the relations between  $S(t)$  and  $\ln S(t)$ :

**Relations for geometric Brownian motion** Suppose the geometric Brownian motion process  $S(t)$  is governed by

$$dS(t) = \mu S(t) dt + \sigma S(t) dz,$$

where  $z$  is a standard Wiener process. Define  $v = \mu - \frac{1}{2}\sigma^2$ . Then  $S(t)$  is lognormal and

$$E\{\ln[S(t)/S(0)]\} = vt$$

$$\text{stddev}\{\ln[S(t)/S(0)]\} = \sigma\sqrt{t}$$

$$E\{S(t)/S(0)\} = e^{\mu t}$$

$$\text{stddev}\{S(t)/S(0)\} = e^{\mu t}(e^{\sigma^2 t} - 1)^{1/2}.$$

## Simulation

A continuous-time price process can be simulated by taking a series of small time periods and then stepping the process forward period by period. There are two natural ways to do this, and they are *not* exactly equivalent.

First, consider the process in standard form defined by equation (13.18). We take a basic period length  $\Delta t$  and set  $S(t_0) = S_0$ , a given initial price at  $t = t_0$ . The corresponding simulation equation is

$$S(t_{k+1}) - S(t_k) = \mu S(t_k) \Delta t + \sigma S(t_k) \epsilon(t_k) \sqrt{\Delta t},$$

where the  $\epsilon(t_k)$ 's are uncorrelated normal random variables of mean 0 and standard deviation 1. This leads to

$$S(t_{k+1}) = \left[ 1 + \mu \Delta t + \sigma \epsilon(t_k) \sqrt{\Delta t} \right] S(t_k), \quad (13.19)$$

which is a multiplicative model, but the random coefficient is normal rather than lognormal, so this simulation method does not produce the lognormal price distributions that are characteristic of the underlying Ito process (in either of its forms). In fact, it is possible for an  $S(t_k)$  to be negative.

A second approach is to use the log (or multiplicative) form (13.15). In discrete form this is

$$\ln S(t_{k+1}) - \ln S(t_k) = v \Delta t + \sigma \epsilon(t_k) \sqrt{\Delta t}.$$

This leads to

$$S(t_{k+1}) = e^{v \Delta t + \sigma \epsilon(t_k) \sqrt{\Delta t}} S(t_k), \quad (13.20)$$

which is also a multiplicative model, but now the random coefficient is lognormal. (In fact the result of (13.20) is distributionally equivalent to the exact result, in the sense that the distribution of  $S(t_{k+1})$  is identical to the distribution that would be attained by the continuous process.)

The two methods are different, but it can be shown that their differences tend to cancel in the long run. Hence in practice, either method is about as good as the other.

**TABLE 13.1**  
**SIMULATION OF PRICE DYNAMICS**

Week	$dz$	$\mu + \sigma dz$	$P_1$	$v + \sigma dz$	$P_2$
0			10.0000		10.0000
1	.06476	.00802	10.0802	.00648	10.0650
2	-.19945	-.00664	10.0132	-.00818	9.9830
3	-.83883	-.04211	9.5916	-.04365	9.5567
4	.49609	.03194	9.8980	.03040	9.8517
5	-.33892	-.01438	9.7557	-.01592	9.6961
6	1.39485	.08180	10.5536	.08026	10.5064
7	.61869	.03874	10.9625	.03720	10.9046
8	.40201	.02672	11.2554	.02518	11.1827
9	-.71118	-.03503	10.8612	-.03656	10.7812
10	.16937	.01382	11.0113	.01228	10.9144
11	1.19678	.07081	11.7910	.06927	11.6973
12	-.14408	-.00357	11.7489	-.00511	11.6377
13	.80590	.04913	12.3261	.04759	12.2049
26	-1.23335	-.06399	13.1428	-.06553	12.9157
39	.68140	.04222	17.6850	.04068	17.3668
52	.69955	.04323	15.1230	.04169	14.7564

The price process is simulated by two methods. Although they differ step by step, the overall results are similar.

**Example 13.3 (Simulation by two methods)** Consider a stock with an initial price of \$10 and having  $v = 15\%$  and  $\sigma = 40\%$ . We take the basic time interval to be 1 week ( $\Delta t = 1/52$ ), and we simulate the stock behavior for 1 year. Both methods described in this subsection were applied using the same random  $\epsilon$ 's, which were generated from a normal distribution of mean 0 and standard deviation 1. Table 13.1 gives the results. The first column shows the random variables  $dz = \epsilon\sqrt{\Delta t}$  for that week. The second column lists the corresponding multiplicative factors. The value  $P_1$  is the simulated price using the standard method as represented by (13.19). The fourth column shows the appropriate exponential factors for the second method (13.20). The value  $P_2$  is the simulated price using that method. Note that even at the first step the results are not identical. However, overall the results are fairly close.

## 13.8 Ito's Lemma\*

We saw that the two Ito equations—for  $S(t)$  and for  $\ln S(t)$ —are different, and that the difference is not exactly what would be expected from the application of ordinary calculus to the transformation of variables from  $S(t)$  to  $\ln S(t)$ ; an additional term  $\frac{1}{2}\sigma^2$  is required. This extra term arises because the random variables have order  $\sqrt{dt}$ , and hence their squares produce first-order, rather than second-order, effects.

There is a systematic method for making such transformations in general, and this is encapsulated in Ito's lemma:

**Ito's lemma** Suppose that the random process  $x$  is defined by the Ito process

$$dx(t) = a(x, t)dt + b(x, t)dz, \quad (13.21)$$

where  $z$  is a standard Wiener process. Suppose also that the process  $y(t)$  is defined by  $y(t) = F(x, t)$ . Then  $y(t)$  satisfies the Ito equation

$$dy(t) = \left( \frac{\partial F}{\partial x} a + \frac{\partial F}{\partial t} + \frac{1}{2} \frac{\partial^2 F}{\partial x^2} b^2 \right) dt + \frac{\partial F}{\partial x} bdz, \quad (13.22)$$

where  $z$  is the same Wiener process as in equation (13.21).

**Proof:** Ordinary calculus would give a formula similar to equation (13.22) but without the term with  $\frac{1}{2}$ .

We shall sketch a rough proof of the full formula. We expand  $y$  with respect to a change  $\Delta y$ . In the expansion we keep terms up to first order in  $\Delta t$ , but since  $\Delta x$  is of order  $\sqrt{\Delta t}$ , this means that we must expand to second order in  $\Delta x$ . We find

$$\begin{aligned} y + \Delta y &= F(x, t) + \frac{\partial F}{\partial x} \Delta x + \frac{\partial F}{\partial t} \Delta t + \frac{1}{2} \frac{\partial^2 F}{\partial x^2} (\Delta x)^2 \\ &= F(x, t) + \frac{\partial F}{\partial x} (a \Delta t + b \Delta z) + \frac{\partial F}{\partial t} \Delta t + \frac{1}{2} \frac{\partial^2 F}{\partial x^2} (a \Delta t + b \Delta z)^2. \end{aligned}$$

The quadratic expression in the last term must be treated in a special way. When expanded, it becomes  $a^2(\Delta t)^2 + 2ab\Delta t \Delta z + b^2(\Delta z)^2$ . The first two terms of this expression are of order higher than 1 in  $\Delta t$ , so they can be dropped. The term  $b^2(\Delta z)^2$  is all that remains. However,  $\Delta z$  has expected value zero and variance  $\Delta t$ , and hence this last term is of order  $\Delta t$  and cannot be dropped. Indeed, it can be shown that, in the limit as  $\Delta t$  goes to zero, the term  $(\Delta z)^2$  is nonstochastic and is equal to  $\Delta t$ . Substitution of this into the previous expansion leads to

$$y + \Delta y = F(x, t) + \left( \frac{\partial F}{\partial x} a + \frac{\partial F}{\partial t} + \frac{1}{2} \frac{\partial^2 F}{\partial x^2} b^2 \right) \Delta t + \frac{\partial F}{\partial x} b \Delta z.$$

Taking the limit and using  $y = F(x, t)$  yields Ito's equation, (13.22). ■

**Example 13.4 (Stock dynamics)** Suppose that  $S(t)$  is governed by the geometric Brownian motion

$$dS = \mu S dt + \sigma S dz.$$

Let us use Ito's lemma to find the equation governing the process  $F(S(t)) = \ln S(t)$ .

We have the identifications  $a = \mu S$  and  $b = \sigma S$ . We also have  $\partial F/\partial S = 1/S$  and  $\partial^2 F/\partial S^2 = -1/S^2$ . Therefore according to equation (13.22),

$$\begin{aligned} d \ln S &= \left( \frac{a}{S} - \frac{1}{2} \frac{b^2}{S^2} \right) dt + \frac{b}{S} dz \\ &= \left( \mu - \frac{1}{2} \sigma^2 \right) dt + \sigma dz, \end{aligned}$$

which agrees with our earlier result.

## 13.9 Binomial Lattice Revisited

Let us consider again the binomial lattice model shown in Figure 13.8 (which is identical to Figure 13.1). The model is analogous to the multiplicative model discussed earlier in this chapter, since at each step the price is multiplied by a random variable. In this case, the random variable takes only the two possible values  $u$  and  $d$ . We can find suitable values for  $u, d$ , and  $p$  by matching the multiplicative model as closely as possible. This is done by matching both the expected value of the logarithm of a price change and the variance of the logarithm of the price change.<sup>6</sup>

To carry out the matching, it is only necessary to ensure that the random variable  $S_1$ , which is the price after the first step, has the correct properties since the process is identical thereafter. Taking  $S(0) = 1$ , we find by direct calculation that

$$\begin{aligned} E(\ln S_1) &= p \ln u + (1-p) \ln d \\ \text{var}(\ln S_1) &= p(\ln u)^2 + (1-p)(\ln d)^2 - [p \ln u + (1-p) \ln d]^2 \\ &= p(1-p)(\ln u - \ln d)^2. \end{aligned}$$

Therefore the appropriate parameter matching equations are

$$pU + (1-p)D = \nu \Delta t \quad (13.23)$$

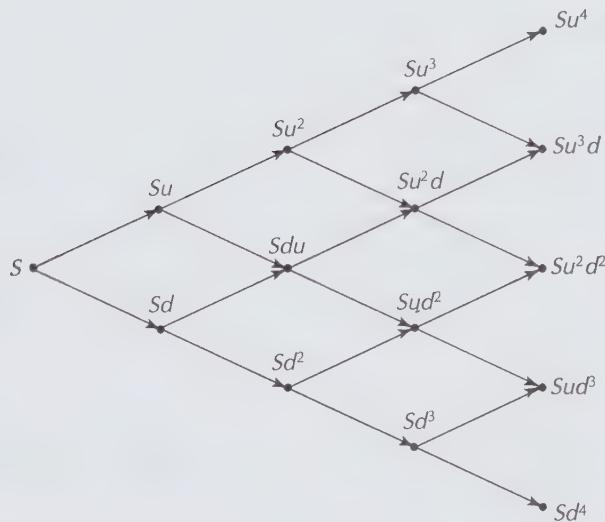
$$p(1-p)(U - D)^2 = \sigma^2 \Delta t \quad (13.24)$$

where  $U = \ln u$  and  $D = \ln d$ .

Notice that three parameters are to be chosen:  $U, D$ , and  $p$ ; but there are only two requirements. Therefore there is one degree of freedom. One way to use this freedom is to set  $D = -U$  (which is equivalent to setting  $d = 1/u$ ). In this case equations

---

<sup>6</sup> For the lattice, the probability of attaining the various end nodes of the lattice is given by the binomial distribution. Specifically, the probability of reaching the value  $Su^k d^{n-k}$  is  $\binom{n}{k} p^k (1-p)^{n-k}$ , where  $\binom{n}{k} = \frac{n!}{(n-k)!k!}$  is the binomial coefficient. This distribution approaches (in a certain sense) a normal distribution for large  $n$ . The logarithm of the final prices is of the form  $k \ln u + (n-k) \ln d$ , which is linear in  $k$ . Hence the distribution of the end point prices can be considered to be nearly lognormal.



**FIGURE 13.8 Binomial lattice stock model.** At each step the stock price  $S$  either increases to  $uS$  or decreases to  $dS$ .

(13.23) and (13.24) reduce to

$$(2p - 1)U = v\Delta t$$

$$4p(1-p)U^2 = \sigma^2\Delta t.$$

If we square the first equation and add it to the second, we obtain

$$U^2 = \sigma^2\Delta t + (v\Delta t)^2.$$

Substituting this in the first equation, we may solve for  $p$  directly. The resulting solutions to the parameter-matching equations are

$$p = \frac{1}{2} + \frac{v\Delta t/2}{\sqrt{\sigma^2\Delta t + (v\Delta t)^2}}$$

$$\ln u = \sqrt{\sigma^2\Delta t + (v\Delta t)^2}$$

$$\ln d = -\sqrt{\sigma^2\Delta t + (v\Delta t)^2}. \quad (13.25)$$

For small  $\Delta t$ , the equation (13.25) can be approximated as

$$p = \frac{1}{2} + \frac{1}{2} \left( \frac{v}{\sigma} \right) \sqrt{\Delta t}$$

$$u = e^{\sigma\sqrt{\Delta t}}$$

$$d = e^{-\sigma\sqrt{\Delta t}}. \quad (13.26)$$

These are the values presented in Section 13.1.

## 13.10 Summary

A simple and versatile model of asset dynamics is the binomial lattice. In this model an asset's price is assumed to be multiplied either by the factor  $u$  or by the factor  $d$ , the choice being made each period according to probabilities  $p$  and  $1 - p$ , respectively. This model is used extensively in theoretical developments and as a basis for computing solutions to investment problems.

Another broad class of models are those where the asset price may take on values from a continuum of possibilities. The simplest model of this type is the additive model. If the random inputs of this model are normal random variables, the asset prices are also normal random variables. This model has the disadvantage, however, that prices may be negative.

A better model is the multiplicative model of the form  $S(k+1) = u(k)S(k)$ . If the multiplicative inputs  $u(k)$  are lognormal, then the future prices  $S(k)$  are also lognormal. The model can be expressed in the alternative form as  $\ln S(k+1) - \ln S(k) = \ln u(k)$ .

By letting the period length tend to zero, the multiplicative model becomes the Ito process  $d \ln S(t) = \nu dt + \sigma dz(t)$ , where  $z$  is a normalized Wiener process. This special form of an Ito process is called geometric Brownian motion. This model can be expressed in the alternative (but equivalent) form  $dS(t) = \mu S(t)dt + \sigma S(t)dz(t)$ , where  $\mu = \nu + \frac{1}{2}\sigma^2$ .

Ito processes are useful representations of asset dynamics. An important tool for transforming such processes is Ito's lemma: If  $x(t)$  satisfies an Ito process, and  $y(t)$  is defined by  $y(t) = F(x, t)$ , Ito's lemma specifies the process satisfied by  $y(t)$ .

A binomial lattice model can be considered to be an approximation to an Ito process. The parameters of the lattice can be chosen so that the mean and standard deviation of the logarithm of the return agree in the two models.

## Exercises

- (Stock lattice) A stock with current value  $S(0) = 100$  has an expected growth rate of its logarithm of  $\nu = 12\%$  and a volatility of that growth rate of  $\sigma = 20\%$ . Find suitable parameters of a binomial lattice representing this stock with a basic elementary period of 3 months. Draw the lattice and enter the node values for 1 year. What are the probabilities of attaining the various final nodes?
- (Time scaling) A stock price  $S$  is governed by the model

$$\ln S(k+1) = \ln S(k) + w(k),$$

where the period length is 1 month. Let  $\nu = E[w(k)]$  and  $\sigma^2 = \text{var}[w(k)]$  for all  $k$ . Now suppose the basic period length is changed to 1 year. Then the model is

$$\ln S(K+1) = \ln S(K) + W(K),$$

where each movement in  $K$  corresponds to 1 year. What is the natural definition of  $W(K)$ ? Show that  $E[W(K)] = 12\nu$  and  $\text{var}[W(K)] = 12\sigma^2$ . Hence parameters scale in proportion to time.

3. (Arithmetic and geometric means) Suppose that  $v_1, v_2, \dots, v_n$  are positive numbers. The *arithmetic mean* and the *geometric mean* of these numbers are, respectively,

$$v_A = \frac{1}{n} \sum_{i=1}^n v_i \quad \text{and} \quad v_G = \left( \prod_{i=1}^n v_i \right)^{1/n}.$$

- (a) It is always true that  $v_A \geq v_G$ . Prove this inequality for  $n = 2$ .  
 (b) If  $r_1, r_2, \dots, r_n$  are rates of return of a stock in each of  $n$  periods, the arithmetic and geometric mean rates of return are likewise

$$r_A = \frac{1}{n} \sum_{i=1}^n r_i \quad \text{and} \quad r_G = \left( \prod_{i=1}^n (1 + r_i) \right)^{1/n} - 1.$$

Suppose \$40 is invested. During the first year it increases to \$60 and during the second year it decreases to \$48. What are the arithmetic and geometric mean rates of return over the 2 years?

- (c) When is it appropriate to use these means to describe investment performance?  
 4. (Complete the square  $\diamond$ ) Suppose that  $u = e^w$ , where  $w$  is normal with expected value  $\bar{w}$  and variance  $\sigma^2$ . Then

$$\bar{u} = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} e^w e^{-(w-\bar{w})^2/2\sigma^2} dw.$$

Show that

$$w - \frac{(w-\bar{w})^2}{2\sigma^2} = -\frac{1}{2\sigma^2}[w - (\bar{w} + \sigma^2)]^2 + \bar{w} + \frac{\sigma^2}{2}.$$

Use the fact that

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} e^{-(x-\bar{x})^2/2\sigma^2} dx = 1$$

to evaluate  $\bar{u}$ .

5. (Log variance  $\diamond$ ) Use the method of Exercise 4 to find the variance of a lognormal variable in terms of the parameters of the underlying normal variable.  
 6. (Expectations) A stock price is governed by geometric Brownian motion with  $\mu = .20$  and  $\sigma = .40$ . The initial price is  $S(0) = 1$ . Evaluate the four quantities

$$\begin{aligned} E[\ln S(1)], & \quad \text{stdev}[\ln S(1)] \\ E[S(1)], & \quad \text{stdev}[S(1)]. \end{aligned}$$

7. (Application of Ito's lemma) A stock price  $S$  is governed by

$$dS = aS dt + bS dz,$$

where  $z$  is a standardized Wiener process. Find the process that governs

$$G(t) = S^{1/2}(t).$$

8. (Reverse check) Gavin Jones was mystified by Ito's lemma when he first studied it, so he tested it. He started with  $S$  governed by

$$dS = \mu S dt + \sigma S dz$$

and found that  $Q = \ln S$  satisfies

$$dQ = (\mu - \frac{1}{2}\sigma^2) dt + \sigma dz.$$

He then applied Ito's lemma to this last equation using the change of variable  $S = e^Q$ . Duplicate his calculations. What did he get?

9. (Forward process) Suppose  $F(S, t)$  is the forward price of a commodity with no storage cost and governed by  $dS(t) = \mu dt + \sigma dz$  and terminating at time  $T$ . What is the process for  $F$ ? (Give the process in terms of  $F$ , not  $S$ .)
10. (Odd process) Let  $w = e^{\sigma z - \frac{1}{2}t^2}$ , where  $z$  is a standard Wiener process. Find the equation governing  $w$ .
11. (Expected Ito) Consider an asset whose price follows the geometric Brownian motion process

$$dS(t) = \mu S(t) dt + \sigma S(t) dz,$$

where  $z$  is a standard Wiener process.

- (a) At time  $t$  (when  $S(t)$  is known), what is the expected value of the asset's price at the future time  $T$ ? This is denoted  $E_t[S(T)]$ .
- (b) Let  $W(t) = E_t[S(T)]$ . What is the process governing  $W(t)$ ?
- (c) Let  $W(0) = 1$ . What are  $E_0[W(t)]$  and  $E_0[\ln W(t)]$ ?
- (d) What are  $\text{Var}_0[W(t)]$  and  $\text{Var}_0[\ln W(t)]$ ?
12. (Fix  $p$ ) An alternative to using  $d = 1/u$  in a binomial model is to use the available degree of freedom by setting  $p = 1/2$ .
- (a) Let  $p = 1/2$ , and find the values of  $u$  and  $d$  that satisfy the matching conditions  $pU + (1-p)D = v\Delta t$  and  $p(1-p)(U-D)^2 = \sigma^2\Delta t$ , where  $U = \ln u$  and  $D = \ln d$ . [For purposes of comparison, use  $\hat{u}$  and  $\hat{d}$  for the resulting alternate values of  $u$  and  $d$ .]
- (b) Which lattice approximation is preferable in applications, and why?
13. (Two simulations ◊) A useful expansion is

$$e^x = 1 + x + \frac{1}{2}x^2 + \dots$$

Use this to express the exponential in equation (13.20) in linear terms of powers of  $\Delta t$  up to first order. Note that this differs from the expression in (13.19), so conclude that the standard form and the multiplicative (or lognormal) form of simulation are different even to first order. Show, however, that the expected values of the two expressions are identical to first order, and hence, over the long run the two methods should produce similar results.

14. (A simulation experiment ⊕) Consider a stock price  $S$  governed by the geometric Brownian motion process

$$\frac{dS}{S(t)} = .10dt + .30dz.$$

- (a) Using  $\Delta t = 1/12$  and  $S(0) = 1$ , simulate several (i.e., *many*) years of this process using either method, and evaluate

$$\frac{1}{t} \ln[S(t)]$$

as a function of  $t$ . Note that it tends to a limit  $p$ . What is the theoretical value of this limit?

- (b) About how large must  $t$  be to obtain two-place accuracy?  
(c) Evaluate

$$\frac{1}{t} [\ln S(t) - pt]^2$$

as a function of  $t$ . Does this tend to a limit? If so, what is its theoretical value?

## References

For a good overview of stock models similar to this chapter, see [1]. For greater detail on stochastic processes see [2], and for general information of how stock prices actually behave, see [3].

There are numerous textbooks on probability theory that discuss the normal distribution and the lognormal distribution. A classic is [4]. The book by Wiener [5] was responsible for inspiring a great deal of serious theoretical and practical work on issues involving Wiener processes. Ito's lemma was first published in [6] and later in [7].

1. Hull, J. C. (2008), *Options, Futures, and Other Derivative Securities*, 7th ed., Prentice Hall, Englewood Cliffs, NJ.
2. Björk, Thomas (2004), *Arbitrage Theory in Continuous Time*, 2nd ed., Oxford University Press, Oxford.
3. Cootner, P. H., Ed. (1964), *The Random Character of Stock Market Prices*, M.I.T. Press, Cambridge, MA.
4. Feller, W. (1950), *Probability Theory and Its Applications*, vols 1 and 2, Wiley, New York.
5. Wiener, N. (1950), *Extrapolation, Interpolation, and Smoothing of Stationary Time Series*, Technology Press, M.I.T., Cambridge, MA, and Wiley, New York.
6. Ito, K. (1951), "On a Formula Concerning Stochastic Differentials," *Nagoya Mathematics Journal*, 3, 55–65.
7. Ito, K. (1961), *Lectures on Stochastic Processes*, Tata Institute of Fundamental Research, India.

## BASIC OPTIONS THEORY

**A**n option is the right, but not the obligation, to buy (or sell) an asset under specified terms. Usually there are a specified price and a specified period of time over which the option is valid. An example is the option to purchase, for a price of \$200,000, a certain house, say, the one you are now renting, anytime within the next year. An option that gives the right to purchase something is called a **call** option, whereas an option that gives the right to sell something is called a **put**. Usually an option itself has a price; frequently we refer to this price as the option **premium**, to distinguish it from the purchase or selling price specified in the terms of the option. The premium may be a small fraction of the price of the optioned asset. For example, you might pay \$15,000 for the option to purchase the house at \$200,000. If the option holder actually does buy or sell the asset according to the terms of the option, the option holder is said to **exercise** the option. The original premium is not recovered in any case.

An option is a derivative security whose underlying asset is the asset that can be bought or sold, such as the house in our example. The ultimate financial value of an option depends on the price of the underlying asset at the time of possible exercise. For example, if the house is worth \$300,000 at the end of the year, the \$200,000 option is then worth \$100,000, because you could buy the house for \$200,000 and immediately sell it for \$300,000 for a profit of \$100,000.

Options have a long history in commerce, since they provide excellent mechanisms for controlling risk, or for locking up resources at a minimal fee. The following story, quoted from Aristotle,<sup>1</sup> is a favorite of professors who write about investments.

---

<sup>1</sup> Aristotle, *Politics*, Book 1, Chapter 11, Jowett translation. Quoted in Gastineau (1975).

There is an anecdote of Thales the Milesian and his financial device, which involves a principle of universal application, but is attributed to him on account of his reputation for wisdom. He was reproached for his poverty, which was supposed to show that philosophy was of no use. According to the story, he knew by his skill in the stars while it was yet winter that there would be a great harvest of olives in the coming year; so, having a little money, he gave deposits for the use of all the olive presses in Chios and Miletus, which he hired at a low price because no one bid against him. When the harvest time came, and many wanted them all at once and of a sudden, he let them out at any rate which he pleased, and made a quantity of money. Thus he showed the world that philosophers can easily be rich if they like ...

Another classic example is associated with the Dutch *tulip mania* in about 1600. Tulips were prized for their beauty, and this led to vigorous speculation and escalation of prices. Put options were used by growers to guarantee a price for their bulbs, and call options were used by dealers to assure future prices. The market was not regulated in any way and finally crashed in 1636, leaving options with a bad reputation.

Options are now available on a wide assortment of financial instruments (such as stocks and bonds) through regulated exchanges. However, options on physical assets are still very important. In addition, there are many implied or hidden options in other financial situations. An example is the option to extract oil from an oil well or leave it in the ground until a better time, or the option to accept a mortgage guarantee or renegotiate. These situations can be fruitfully analyzed using the theory of options explained in this chapter.

## 14.1 Option Concepts

The specifications of an option include, first, a clear description of what can be bought (for a call) or sold (for a put). For options on stock, each option is usually for 100 shares of a specified stock. Thus a call option on IBM is the option to buy 100 shares of IBM. Second, the exercise price, or **strike price**, must be specified. This is the price at which the asset can be purchased upon exercise of the option. For IBM stock the exercise price might be \$70, which means that each share can be bought at \$70. Third, the period of time for which the option is valid must be specified—defined by the expiration date. Hence an option may be valid for a day, a week, or several months. There are two primary conventions regarding acceptable exercise dates before expiration. An **American option** allows exercise at any time before and including the expiration date. A **European option** allows exercise only on the expiration date. The terms *American* and *European* refer to the different ways most stock options are structured in America and in Europe, but the words have become standard for the two different types of structures, no matter where they are issued. There are some

European-style options in America. For example, if the option to buy a house in one year states that the sale must be made in exactly one year and not sooner, the house option can be referred to as a European option.

These four features—the description of the asset, whether a call or a put, the exercise price, and the expiration date (including whether American or European in style)—specify the details of an option. A final, but somewhat separate, feature is the price of the option itself—the premium. If an option is individually tailored, this premium price is established as part of the original negotiation and is part of the contract. If the option is traded on an exchange, the premium is established by the market through supply and demand, and this premium will vary according to trading activity.

There are two sides to any option: the party that grants the option is said to **write** an option, whereas the party that obtains the option is said to purchase it. The party purchasing an option faces no risk of loss other than the original purchase premium. However, the party that writes the option may face a large loss, since this party must buy or sell this asset at the specified terms if the option is exercised. In the case of an exercised call option, if the writer does not already own the asset, he must purchase it in order to deliver it at the specified strike price, which may be much lower than the current market price. Likewise, in the case of an exercised put option, the writer must accept the asset for the strike price, which could be much higher than the current market price.

Options on many stocks are traded on an exchange. In this case individual option trades are made through a broker who trades on the exchange. The exchange clearing-house guarantees the performance of all parties. Because of the risk associated with options, an option writer is required to post **margin** (a security deposit) guaranteeing performance.<sup>2</sup>

Exchange-traded options are listed in the financial press and online. A listing of options for Apple Inc. is shown in Figure 14.1. Several different options are shown; some are puts and some are calls, and they have a variety of strike prices and expiration dates. The expiration dates are September, October, and January. The date of expiration follows an expiration calendar. The September options listed for Apple have an expiration date of September 23. All prices are quoted on a per-share basis, although a single option contract is for 100 shares. This particular September was one of high volatility in the market.

As with futures contracts, options on financial securities are rarely exercised, with the underlying security being bought or sold. Instead, if the price of the security moves in a favorable direction, the option price (the premium) will increase accordingly, and most option holders will elect to sell their options before maturity.

There are many details with regard to options trading, governing special situations such as stock splits, dividends, position limits, and specific margin requirements. These must be checked before engaging in serious trading of options.

---

<sup>2</sup> The initial margin level is often 50% of the stock value of the option, with a maintenance level of 25%.

Close	Strike	September		October		January	
		Vol.	Close	Vol.	Close	Vol.	Close
413.35	p 350	3.7k	1.67	1.1k	4.30	592	11.00
413.35	p 360	5.7k	2.42	562	5.64	512	13.10
413.35	p 370	3.8k	3.40	1.1k	7.20	340	15.70
413.35	c 380	6.4k	38.60	492	43.60	1.5k	52.50
413.35	p 380	3.6k	4.95	853	9.60	311	18.65
413.35	c 390	4.2k	30.75	1.0k	36.09	1.3k	46.20
413.35	p 390	3.3k	7.10	973	12.25	451	21.75
413.35	p 400	6.1k	10.10	1.0k	15.85	687	26.00
413.35	c 400	14k	22.10	4.1k	29.80	4.8k	39.95
413.35	c 410	1.1k	17.75	2.4k	24	1.1k	34.30
413.35	p 410	6.5k	13.95	1.1k	20	281	21.95
413.35	c 420	21k	12.80	4.4k	18.60	1.8k	28.95
413.35	p 420	4.9k	19.95	1.4k	25	369	33.85

**FIGURE 14.1 Options on Apple Inc., September 21, 2011.** Puts and calls for various strike prices and maturity dates are listed. For each option the daily volume and closing price are shown. The closing price of Apple is shown in the first column.

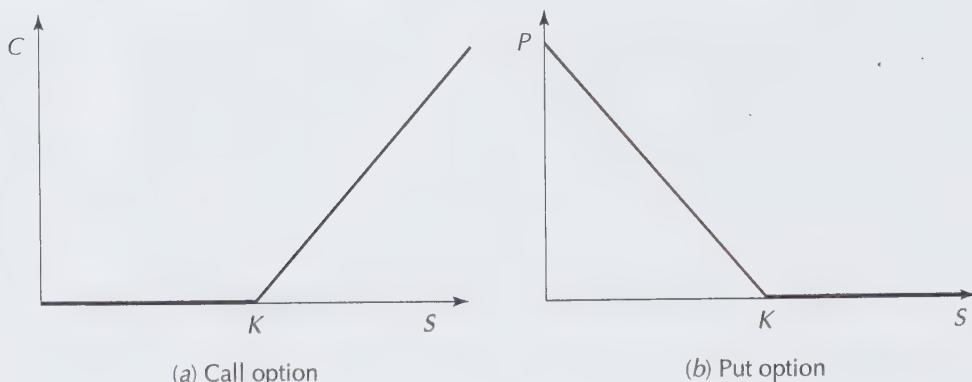
Source: *Investor's Business Daily*.

However, the present overview is sufficient for understanding the basic mechanics of options.

## 14.2 The Nature of Option Values

A primary objective of this chapter is to show how to determine the value of an option on a financial security. Such a determination is a fascinating and creative application of the fundamental principles that we have studied so far. Hence options theory is important *partly* because options themselves are important financial instruments, but also partly because options theory shows how the fundamental principles of investment science can be taken to a new level—a level where dynamic structure is fundamental. In this section we examine in a qualitative manner the nature of option prices. This will prepare us for the deeper analysis that follows in subsequent sections.

Suppose that you own a call option on a stock with a strike price of  $K$ . Suppose that on the expiration date the price of the underlying stock is  $S$ . What is the value of the option at that time? It is easy to see that if  $S < K$ , then the option value is zero. This is because under the terms of the option, you could exercise the option and purchase the stock for  $K$ , but by not exercising the option you could buy the stock on the open market for the lower price of  $S$ . Hence you would not exercise the option. The option is worthless. On the other hand, if  $S > K$ , then the option does have value. By exercising the option you could buy the stock at a price  $K$  and then sell that stock on the market for the larger price  $S$ . Your profit would be  $S - K$ , which is therefore the value of the option. We handle both cases together by writing the value of the call



**FIGURE 14.2 Value of option at expiration.** A call has value if  $S > K$ . A put has value if  $S < K$ .

at expiration as

$$C = \max(0, S - K), \quad (14.1)$$

which means that  $C$  is equal to the maximum of the values 0 or  $S - K$ . We therefore have an explicit formula for the value of a call option at expiration as a function of the price of the underlying security  $S$ . This function is shown in Figure 14.2(a). The figure shows that for  $S < K$ , the value is zero, but for  $S > K$ , the value of the option increases linearly with the price, on a one-for-one basis.

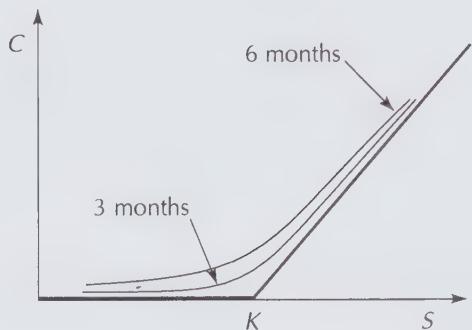
The result is reversed for a put option. A put option gives one the right, but not the obligation, to sell an asset at a given strike price. Suppose you own a put option on a stock with a strike price of  $K$ . In this case if the price  $S$  of the stock at expiration satisfies  $S > K$ , then this option is worthless. By exercising the option you could sell the stock for a price  $K$ , whereas in the open market you could sell the stock for the greater price  $S$ . Hence you would not exercise the option. On the other hand, if the price of the stock is less than the strike price, the put option does have value. You could buy the stock on the market for a price  $S$  and then exercise the option to sell that same stock for a greater price  $K$ . Your profit would be  $K - S$ , which is therefore the value of the option. The general formula for the value of a put at expiration is

$$P = \max(0, K - S). \quad (14.2)$$

This function is illustrated in Figure 14.2(b). Note that the value of a put is bounded, whereas the payoff of a call is unbounded. Conversely, when writing a call, the potential for *loss* is unbounded.

We say that a call option is **in the money**, **at the money**, or **out of the money**, depending on whether  $S > K$ ,  $S = K$ , or  $S < K$ , respectively. The terminology applies at any time; but at expiration the terms describe the nature of the option value. Puts have the reverse terminology, since the payoffs at exercise are positive if  $S < K$ .

**FIGURE 14.3 Option price curve with various times to expiration.**  
At a given stock price  $S$ , the value of a call option increases as the time to expiration increases.



## Time Value of Options

The preceding analysis focused on the value of an option at expiration. This value is derived from the basic structure of an option. However, even European options (which cannot be exercised except at expiration) have value at earlier times, since they provide the potential for future exercise. As an example in a highly volatile environment, consider the September call on Apple Inc. with strike price of 420. This option, which expires in just 2 days, is out of the money, yet it has a price of 12.80. (The next day Apple closed at 412.14, and the value of this call closed at 6.90.)

Generally, when there is a positive time to expiration, the value of a call option as a function of the stock price is a smooth curve rather than the decidedly kinked curve that applies at expiration. This smooth curve can be determined by estimation, using data of actual option prices. Such estimation shows that the option price curve for any given expiration period looks something like the curves shown in Figure 14.3. In this figure the heavy kinked line represents the value of a call at expiration. The higher curves correspond to different times to expiration. The first curve is for a call with 3 months to expiration, whereas the next higher one is for 6 months. The curves get higher with increasing length to expiration, since additional time provides a greater chance for the stock to rise in value, increasing the final payoff. However, the effect of additional time is diminished when the stock price is either much smaller or much greater than the strike price  $K$ . When the stock price  $S$  is much lower than  $K$ , there is little chance that  $S$  will rise above  $K$ , so the option value remains close to zero. When  $S$  is much greater than  $K$ , there is little advantage in owning the option over owning the stock itself.

A major objective of this chapter is to determine a theory for option prices. This theory will imply a specific set of curves, such as the ones shown in Figure 14.3.

## Other Factors Affecting the Value of Options

The volatility of the underlying stock is another factor that influences the value of an option significantly. To see this, imagine that you own similar options on two different stocks. Suppose the prices of the two stocks are both \$90, the options have

strike prices of \$100, and there are 3 months to expiration. Suppose, however, that one of these stocks is very volatile and the other is quite placid. Which option has more value? It is clear that the stock with the high volatility has the greatest chance of rising above \$100 in the short period remaining to expiration, and hence its option is the more valuable of the two. We expect therefore that the value of a call option increases with volatility, and we shall verify this in our theoretical development.

What other factors might influence the value of an option? One is the prevailing interest rate (or term structure pattern). Purchasing a call option is in some way a method of purchasing the stock at a reduced price. Hence one saves interest expense. We expect therefore that option prices depend on interest rates.

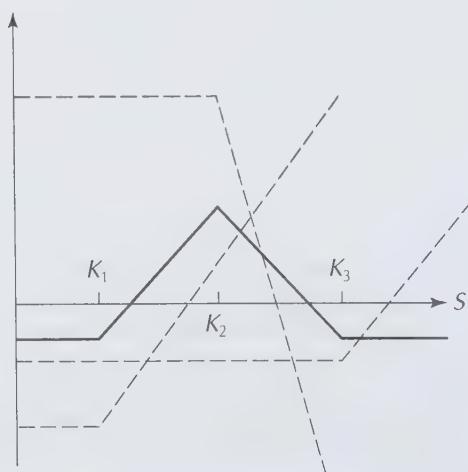
Another factor that would seem to be important is the growth rate of the stock. It seems plausible that higher values of growth would imply larger values for the option. However, perhaps surprisingly, the growth rate does *not* influence the theoretical value of an option. The reason for this will become clear when the theoretical formula is developed.

## 14.3 Option Combinations and Put–Call Parity

It is common to invest in combinations of options in order to implement special hedging or speculative strategies. The payoff curve of such a combination may have any number of connected straight-line segments. This overall payoff curve is formed by combining the payoff functions defined by calls, puts, and the underlying stock itself. The process is best illustrated by an example and a corresponding graph.

**Example 14.1 (A butterfly spread)** One of the most interesting combinations of options is the butterfly spread. It is illustrated in Figure 14.4. The spread is constructed by buying two calls, one with strike price  $K_1$  and another with strike price  $K_3$ , and by selling two units of a call with strike price  $K_2$ , where  $K_1 < K_2 < K_3$ . Usually  $K_2$

**FIGURE 14.4 Profit of butterfly spread.** This spread is formed by buying calls with strike prices  $K_1$  and  $K_3$  and writing two units of a call at  $K_2$ . This combination is useful if one believes that the underlying stock price will stay in a region near  $K_2$ .



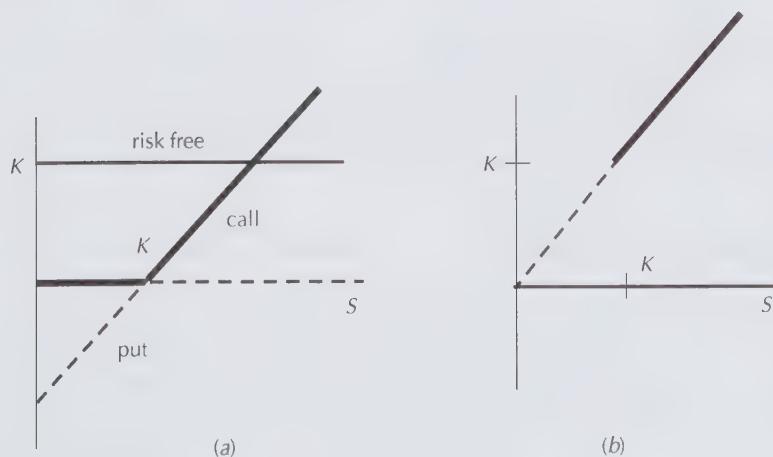
is chosen to be near the current stock price. The figure shows with dashed lines the *profit* (including the payoff and original cost) associated with each of the components. The overall profit function of the combination is the sum of the individual component functions. This particular combination yields a positive profit if the stock price at expiration is close to  $K_2$ ; otherwise the loss is quite small. The payoff of this spread is obtained by lifting the curve up so that the horizontal portions touch the axis, the displacement distance corresponding to the net cost of the options.

The main point here is that by forming combinations of options and stock it is possible to approximate virtually any payoff function by a sequence of straight-line segments. The cost of such a payoff is then just the sum of the costs of the individual components.

## Put-Call Parity

For European options there is a simple theoretical relationship between the prices of corresponding puts and calls. The relationship is found by noting that a combination of a put, a call, and a risk-free loan has a payoff identical to that of the underlying stock.

The combination can be easily imagined: buy one call, sell one put, and lend an amount  $dK$ , where  $d$  is the discount factor for the period. The combination of the first two has a payoff that is a straight line at  $45^\circ$ , passing through  $K$  on the horizontal axis. By lending  $dK$ , we obtain an additional payoff of  $K$ , which lifts the payoff line up so that it is now a  $45^\circ$  line originating at the origin. See Figure 14.5. This final payoff is exactly that of the stock itself, so it must have the value  $S$  of the stock.



**FIGURE 14.5** (a) The payoffs of a call, a short put, and a risk-free payoff of  $K$ . The total initial cost is  $C - P + dK$ . (b) The diagrams for the call and the put are lifted by an amount  $K$  to get the composite payoff, which is the same as the payoff of the underlying stock itself.

**TABLE 14.1**  
**VERIFICATION OF PUT-CALL PARITY**

Expiration	$C - P + K$	
September	$38.60 - 4.95 + 380$	= 413.65
October	$43.60 - 9.60 + 380$	= 414.00
January	$52.50 - 18.65 + 380$	= 413.85

In other words,

$$C - P + dK = S.$$

**Put-call parity** Let  $C$  and  $P$  be the prices of a European call and a European put, both with a strike price of  $K$  and both defined on the same stock with price  $S$ . Put-call parity states that

$$C - P + dK = S,$$

where  $d$  is the discount factor to the expiration date.

**Example 14.2 (Parity check)** Consider the Apple options of Figure 14.1. The put-call parity relation depends on the interest rate, but at that time the short-term rate was essentially zero. Thus we expect that the relation should be

$$C - P + K = S.$$

Let us focus on the put and calls with strike price 380 and evaluate the left-hand side of the parity equation. It should turn out to be the stock price of 413.35. Table 14.1 shows that for all of the three expiration dates, the result is extremely close.

In practice there may be slight mismatches in the parity relation. There are several possible explanations for this. One of the most important is that the stock quotes and option quotes do not come from the same sources. The stock price is the closing price on the stock exchange, whereas the option prices are from the last traded options on the options exchanges; the last trades can occur at different times. Dividends also can influence the parity relation, as discussed in Exercise 2.

## 14.4 Early Exercise

An American option offers the possibility of early exercise, that is, exercise before the expiration date of the option. We prove in this section that for call options on a stock that pays no dividends prior to expiration, early exercise is never optimal, provided that prices are such that no arbitrage is possible.

The result can be seen intuitively as follows. Suppose that we are holding a call option at time  $t$  and expiration is at time  $T > t$ . If the current stock price  $S(t)$  is less than the strike price  $K$ , we would not exercise the option, since we would lose money. If, on the other hand, the stock price is greater than  $K$ , we might be tempted to exercise. However, if we do so we will have to pay  $K$  now to obtain the stock. If we hold the option a little longer and then exercise, we will still obtain the stock for a

price of  $K$ , but we will have earned additional interest on the exercise money  $K$ —in fact, if the stock declines below  $K$  in this waiting period, we will not exercise and be happy that we did not do so earlier.

## 14.5 Single-Period Binomial Options Theory

We now turn to the issue of calculating the theoretical value of an option—an area of work that is called **options pricing theory**. There are several approaches to this problem, based on different assumptions about the market, about the dynamics of stock price behavior, and about individual preferences. The most important theories are based on the no arbitrage principle, which can be applied when the dynamics of the underlying stock take certain forms. The simplest of these theories is based on the binomial model of stock price fluctuations discussed in Chapter 13. This theory is widely used in practice because of its simplicity and ease of calculation. It is a beautiful culmination of the principles discussed in previous chapters.

The basic theory of binomial options pricing has been hinted at in our earlier discussions. We shall develop it here in a self-contained manner, but the reader should notice the connections to earlier sections.

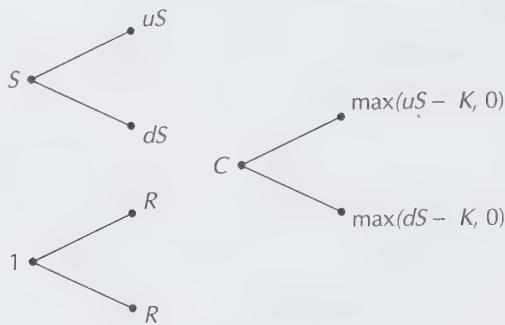
We shall first develop the theory for the single-period case. A single step of a binomial process is all that is used. Accordingly, we suppose that the initial price of a stock is  $S$ . At the end of the period the price will either be  $uS$  with probability  $p$  or  $dS$  with probability  $1 - p$ . We assume  $u > d > 0$ . Also at every period it is possible to borrow or lend at a common risk-free interest rate  $r$ . We let  $R = 1 + r$ . To avoid arbitrage opportunities, we must have

$$u > R > d.$$

To see this, suppose  $R \geq u > d$  and  $0 < p < 1$ . Then the stock performs worse than the risk-free asset, even in the “up” branch of the lattice. Hence one could short \$1.00 of the stock and loan the proceeds, thereby obtaining a profit of either  $R - u$  or  $R - d$ , depending on the outcome state. The initial cost is zero, but in either case the profit is positive, which is not possible if there are no arbitrage opportunities. A similar argument rules out  $u > d \geq R$ .

Now suppose also that there is a call option on this stock with exercise price  $K$  and expiration at the end of the period. To find the value of the call option, we use a no-arbitrage argument by referring to Figure 14.6. This figure shows the binomial lattices for the stock price, the value of a risk-free asset, and the value of the option. All three of these lattices have common arcs, in the sense that all move together along the same arcs. If the stock price moves along the upward arc, then the risk-free asset and the call option both move along their upward arcs as well. The risk-free value is deterministic, but this is treated as if it were a (degenerate) derivative of the stock by just making the value at the end of each arc the same.

Assuming that we know the stock price  $S$ , then all values of these one-step lattices are known except the value of the call  $C$ . This value will be determined from the other values.



**FIGURE 14.6 Three related lattices.** The stock price, the value of a risk-free loan, and the value of a call option all move together on a common lattice, represented here as three separate lattices.

The insight that we use is to note that each of the lattices on the left has only two possible outcomes. By combining various proportions of these two lattices, we can construct any other pattern of outcomes. In particular, we can construct the pattern corresponding to the outcomes of the option.

Let us denote

$$C_u = \max(uS - K, 0) \quad (14.3)$$

$$C_d = \max(dS - K, 0). \quad (14.4)$$

To duplicate these two outcomes, let us purchase  $x$  dollars worth of stock and  $b$  dollars worth of the risk-free asset. At the next time period, this portfolio will be worth either  $ux + Rb$  or  $dx + Rb$ , depending on which path is taken. To match the option outcomes we therefore require

$$ux + Rb = C_u \quad (14.5a)$$

$$dx + Rb = C_d. \quad (14.5b)$$

To solve these equations we subtract the second from the first, obtaining

$$x = \frac{C_u - C_d}{u - d}.$$

From this we easily find

$$b = \frac{C_u - ux}{R} = \frac{uC_d - dC_u}{R(u - d)}.$$

Combining these we find that the value of the portfolio is

$$\begin{aligned} x + b &= \frac{C_u - C_d}{u - d} + \frac{uC_d - dC_u}{R(u - d)} \\ &= \frac{1}{R} \left( \frac{R - d}{u - d} C_u + \frac{u - R}{u - d} C_d \right). \end{aligned}$$

We now use the comparison principle (or, equivalently, the no-arbitrage principle) to assert that the value  $x + b$  must be the value of the call option  $C$ . The reason is that the portfolio we constructed produces exactly the same outcomes as the call option. If the cost of this portfolio were less than the price of the call, we would never purchase the call. Indeed, we could make arbitrage profits by buying this portfolio and selling the call for an immediate gain and no future consequence. If the prices were unequal in the reverse direction, we could just reverse the argument. We conclude therefore that the price of the call is

$$C = \frac{1}{R} \left( \frac{R-d}{u-d} C_u + \frac{u-R}{u-d} C_d \right). \quad (14.6)$$

The portfolio made up of the stock and the risk-free asset that duplicates the outcome of the option is often referred to as a **replicating portfolio**. It replicates the option. This replicating idea can be used to find the value of any security defined on the same lattice; that is, any security that is a derivative of the stock.

There is a simplified way to view equation (14.6). We define the quantity

$$q = \frac{R-d}{u-d}. \quad (14.7)$$

From the relation  $u > R > d$  assumed earlier, it follows that  $0 < q < 1$ . Hence  $q$  can be considered to be a probability. Note also that the coefficients of  $C_u$  and  $C_d$  sum to 1. Hence, equation (14.6) can be written as follows:

**Option pricing formula** *The value of a one-period call option on a stock governed by a binomial lattice is*

$$C = \frac{1}{R} [qC_u + (1-q)C_d]. \quad (14.8)$$

Note that equation (14.8) can be interpreted as stating that  $C$  is found by taking the expected value of the option using the probability  $q$ , and then discounting this value according to the risk-free rate. The probability  $q$  is therefore a **risk-neutral probability**. This procedure of valuation works for all securities. In fact  $q$  can be calculated by making sure that the risk-neutral formula holds for the underlying stock itself; that is, we want

$$S = \frac{1}{R} [quS + (1-q)dS].$$

Solving this equation gives (14.7).

As a suggestive notation, we write equation (14.8) as

$$C(T-1) = \frac{1}{R} \hat{E}[C(T)].$$

Here  $C(T)$  and  $C(T-1)$  are the call values at  $T$  and  $T-1$ , respectively, and  $\hat{E}$  denotes expectation with respect to the risk-neutral probabilities  $q$  and  $1-q$ .

An important, and perhaps initially surprising, feature of pricing formula (14.6) is that it is *independent* of the probability  $p$  of an upward move in the lattice. This is because no trade-off among probabilistic events is made. The value is found by perfectly matching the outcomes of the option with a combination of stock and the risk-free asset. Probability never enters this matching calculation.

This derivation of the option pricing formula is really a special case of the risk-neutral pricing concept discussed in Chapter 11. The pricing formula can be derived very quickly using the concept of linear pricing. By linearity, we can write immediately that

$$C = \frac{a}{R} C_u + \frac{b}{R} C_d$$

for some  $a$  and  $b$ . To avoid arbitrage there must hold  $a \geq 0$ ,  $b \geq 0$ . Applying this to the risk-free asset gives  $1 = a + b$ . Thus the pricing formula is of the form

$$C = \frac{1}{R} \{qC_u + (1 - q)C_d\}$$

for some  $q$  with  $0 \leq q \leq 1$ . Applying this to the stock, we get

$$1 = \frac{1}{R} \{qu + (1 - q)d\},$$

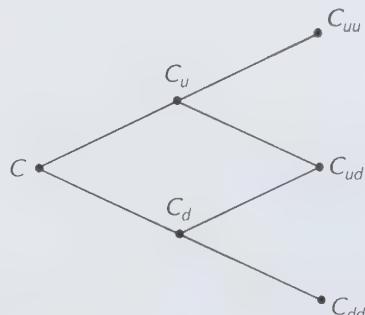
and hence  $q = (R - d)/(u - d)$ .

## 14.6 Multiperiod Options

We now extend the solution method to multiperiod options by working backward one step at a time.

A two-stage lattice representing a two-period call option is shown in Figure 14.7. It is assumed as before that the initial price of the stock is  $S$ , and this price is modified by the up and down factors  $u$  and  $d$  while moving through the lattice. The values shown in the lattice are those of the corresponding call option with strike price  $K$  and expiration time corresponding to the final point in the lattice. The value of the option

**FIGURE 14.7 Two-period option.** The value is found by working backward a step at a time.



is known at the final nodes of the lattice. In particular,

$$C_{uu} = \max(u^2 S - K, 0) \quad (14.9a)$$

$$C_{ud} = \max(udS - K, 0) \quad (14.9b)$$

$$C_{dd} = \max(d^2 S - K, 0). \quad (14.9c)$$

We again define the risk-neutral probability as

$$q = \frac{R - d}{u - d},$$

where  $R$  is the one-period return on the risk-free asset. Then, assuming that we do not exercise the option early (which we already know is optimal, but will demonstrate again shortly), we can find the values of  $C_u$  and  $C_d$  from the single-period calculation given earlier. Specifically,

$$C_u = \frac{1}{R}[qC_{uu} + (1 - q)C_{ud}] \quad (14.10)$$

$$C_d = \frac{1}{R}[qC_{ud} + (1 - q)C_{dd}]. \quad (14.11)$$

Then we find  $C$  by another application of the same risk-neutral discounting formula. Hence,

$$C = \frac{1}{R}[qC_u + (1 - q)C_d].$$

For a lattice with more periods, a similar procedure is used. The single-period, risk-free discounting is just repeated at every node of the lattice, starting from the final period and working backward toward the initial time.

**Example 14.3 (A 5-month call)** Consider a stock with a volatility of its logarithm of  $\sigma = .20$ . The current price of the stock is \$62. The stock pays no dividends. A certain call option on this stock has an expiration date 5 months from now and a strike price of \$60. The current rate of interest is 10%, compounded monthly. We wish to determine the theoretical price of this call using the binomial option approach.

First we must determine the parameters for the binomial model of the stock price fluctuations. We shall take the period length to be 1 month, which means  $\Delta t = 1/12$ . The parameters are found from equations (13.1) to be

$$u = e^{\sigma\sqrt{\Delta t}} = 1.05943$$

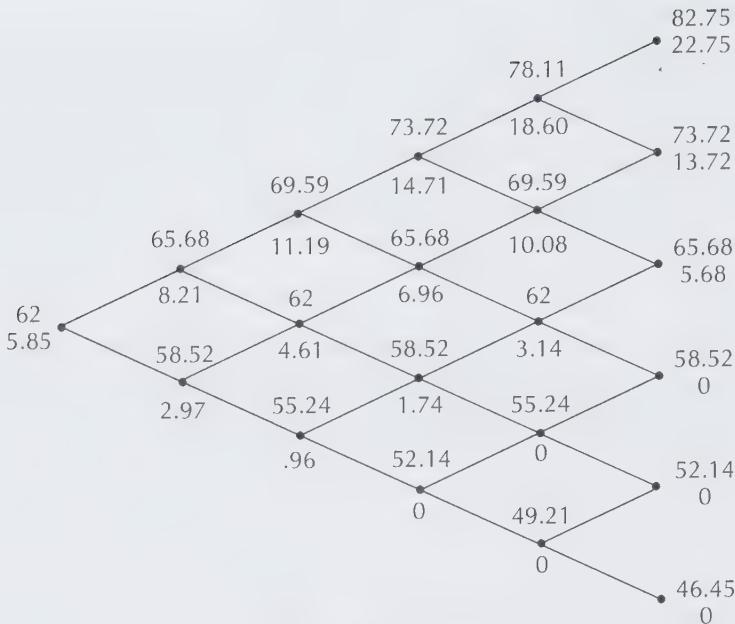
$$d = e^{-\sigma\sqrt{\Delta t}} = .94390$$

$$R = 1 + .1/12 = 1.00833.$$

Then the risk-neutral probability is

$$q = (R - d)/(u - d) = .55770.$$

We now form the binomial lattice corresponding to the stock price at the beginning of each of six successive months (including the current month). This lattice is



**FIGURE 14.8 5-month call using a binomial lattice.** The upper numbers are the stock prices, the lower numbers the option values. The option values are found by working backward through the lattice.

shown in Figure 14.8, with the number above a node being the stock price at that node. Note that an up followed by a down always yields a net multiple of 1.

Next we calculate the call option price. We start at the final time and enter the expiration values of the call below the final nodes. This is the maximum of 0 and  $S - K$ . For example, the entry for the top node is  $82.75 - 60 = 22.75$ .

The values for the previous time are found by the single-step pricing relation. The value of any node at this time is the discounted expected value of two successive values at the next time. The expected value is calculated using the risk-neutral probabilities  $q$  and  $1 - q$ . For example, the value at the top node is  $[.5577 \times 22.75 + (1 - .5577) \times 13.72]/1.00833 = 18.60$ .

We work toward the left, one period at a time, until finally the initial value is reached. In this case we conclude that the price of the option computed this way is \$5.85.

Note that the entire process is independent of the expected growth rate of the stock. This value only enters the binomial model of the stock through the probability  $p$ ; but this probability is not used in the option calculation. Instead it is the risk-neutral probability  $q$  that is used. Note, however, that this independence results from using the small  $\Delta t$  approximation for parameter matching. And indeed, in practice this approximation is almost invariably used (even for  $\Delta t$  equal to 1 year). If the more general matching formula were used, the growth rate would (slightly) influence the result.

## No Early Exercise\*

In the preceding example we assumed (rightly) that the option would never be exercised early. We can prove this directly from the binomial equations under the assumption that  $R > 1$ . From the basic payoff structure we see that

$$C_{uu} \geq u^2 S - K$$

$$C_{ud} \geq udS - K$$

$$C_{dd} \geq d^2 S - K.$$

Hence,

$$\begin{aligned} C_u &\geq [u^2 qS + ud(1-q)S - K]/R \\ &= u[qu + (1-q)d]S/R - K/R \\ &> uS - K. \end{aligned}$$

Likewise,

$$C_d > dS - K.$$

If the option were exercised at the end of the first period of the two-period lattice shown in Figure 14.6, we would obtain  $uS - K$  or  $dS - K$ , depending on which node was active at the time. These inequalities show that the value of the option at the end of one period is greater than the amount that would be obtained by exercise at that period. Hence we should not exercise the option.

If the lattice had more periods, these inequalities would extend to the next forward period as well. Hence, in general, by an inductive process it can be shown that it is never optimal to exercise the option.

The argument against early exercise does not hold for all options; in some cases an additional operation must be incorporated in the recursive process of value calculation. This is explained in the next section.

## 14.7 More General Binomial Problems

The binomial lattice method for calculating the value of an option is extremely simple and highly versatile. For this reason it has become a common tool in the investment and financial community. The method is simplest when applied to a call option on a non-dividend-paying stock, as illustrated in the previous section. This section shows how the basic method can be extended to more complex situations.

### Put Options

The method for calculating the values of European put options is analogous to that for call options. The main difference is that the terminal values for the option are different. But once these are specified, the recursive procedure works in a similar way.

**FIGURE 14.9 Calculation of a 5-month put option price.** The put values in the lower portion of the figure are found by working backward. Boldface entries indicate points where it is optimal to exercise the option.

	62.00	65.68	69.59	73.72	78.11	82.75	
		58.52	62.00	65.68	69.59	73.72	
			55.24	58.52	62.00	65.68	
				Stock price	52.14	55.24	58.52
						49.21	52.14
							46.45
	1.56	0.61	0.12	0.00	0.00	0.00	
		2.79	1.23	0.28	0.00	0.00	
			4.80	2.45	0.65	0.00	
				Put option	<b>7.86</b>	<b>4.76</b>	<b>1.48</b>
					<b>10.79</b>	<b>7.86</b>	
						<b>13.55</b>	

For an American put, early exercise may be optimal. This is easily accounted for in the recursive process as follows: At each node, first calculate the value of the put using the discounted risk-neutral formula; then calculate the value that would be obtained by immediate exercise of the put; finally, select the larger of these two values as the value of the put at that node.

**Example 14.4 (A 5-month put)** We consider the same stock that was used to evaluate the 5-month call option of Example 14.3, but now we evaluate a 5-month American put option with a strike price of  $K=\$60$ . Recall that the critical parameters were  $R = 1.008333$ ,  $q = .55770$ ,  $u = 1.05943$ , and  $d = .94390$ . Binomial lattice calculations can be very conveniently carried out with a spreadsheet program. Hence we often show lattices in spreadsheet form rather than as graphical diagrams. This allows us to show larger lattices in a restricted space, and it also indicates more directly how calculations are organized.

The binomial lattice for the stock price is shown in the top portion of Figure 14.9. In this figure an up move is made by moving directly to the right, and a down move is made by moving to the right and down one step.

To calculate the value of the put option, we again work backward, constructing a new lattice below the stock price lattice. The final values (those of the last column) are, in this case, the maximum of 0 and  $K - S$ . We then work toward the left, one column at a time. To find the value of an element we first calculate the discounted expected value as before, using the risk-neutral probabilities. Now, however, we must also check whether this value would be exceeded by  $K - S$ , which is what could be obtained by exercising the option at the current point. We assign the larger of the two values to this current node. For example, consider the fourth entry in the second to last column. The discounted expected value there is  $[.5577 \times 1.48 + (1 - .5577) \times 7.86]/1.00833 = 4.266$ . The exercise value is  $60 - 55.24 = 4.76$ . The larger of these is 4.76, and that is what is entered in the value lattice. If the larger value is obtained by exercising, we may also wish to indicate this on the lattice, which in our figure is done by using boldface for the entries corresponding to exercise points. (Alternatively, a separate

lattice consisting of 0's and 1's can be constructed to indicate the exercise points.) In our example we see that there are several points at which exercise is optimal. The value of the put is the first entry of the lattice, namely, \$1.56.

Intuitively, early exercise of a put may be optimal because the upside profit is bounded. Clearly, for example, if the stock price falls to zero, one should exercise there, since no greater profit can be achieved. A continuity argument can be used to infer that it is optimal to exercise if the stock price gets close to zero.

## Dividend and Term Structure Problems\*

Many other problems can be treated with the binomial lattice model by allowing the parameters of the model to vary from node to node. This does not change the basic structure of the computational method. It merely means that the risk-neutral probabilities and the discount factor may differ from period to period.

One example is the evaluation of a call option on a stock that pays a dividend. If the dividend is proportional to the value of the stock—say, the dividend is  $\delta S$  and is paid at time  $k$ —then in the stock price lattice we just change the factors  $u$  and  $d$  for the period ending at  $k$  to  $u(1 - \delta)$  and  $d(1 - \delta)$ . If the dividend is known in advance to be a fixed amount  $D$ , then this technique will not work directly, but the lattice approach can still be used. (See Exercise 5.)

The parameters also vary when the interest rate is not constant. In this case the appropriate single-period rate for a given period (the implied short rate) should be used. This will change the value of  $R$  and hence also the value of  $q$ .

## Futures Options\*

Are we ready to consider a futures option—that is, an option on a futures contract? This may at first sound complicated; but we shall find that futures options are quite simple to analyze, and study of the analysis should help develop a fuller understanding of the risk-neutral pricing process. The best way to study the analysis is to consider an example.

**Example 14.5 (A futures contract)** Suppose that a certain commodity (which can be stored without cost and is in ample supply) has a current price of \$100, and the price process is described by a monthly binomial lattice with parameters  $u = 1.02$ ,  $d = .99$ , and  $R = 1.01$ . The actual probabilities are not important for our analysis. This lattice, for 6 months into the future, is shown in the upper left-hand corner of Figure 14.9. We can immediately calculate the risk-neutral probabilities to be  $q = (R - d)/(u - d) = \frac{2}{3}$  and  $1 - q = \frac{1}{3}$ .

Let us compute the lattice of the corresponding futures prices for a futures contract that expires in the sixth month. This lattice is shown in the lower left-hand side of Figure 14.10. One way to compute this lattice is to use the result of Chapter 12 that the futures price is equal to the current commodity price amplified by interest rate growth over the remaining period of the contract. Hence the futures price at time

0	1	2	3	4	5	6	0	1	2	3	4	5	6
100.00	102.00	104.04	106.12	108.24	110.41	112.62	4.16	5.05	6.04	7.12	8.25	9.42	10.62
	99.00	100.98	103.00	105.06	107.16	109.30		2.50	3.21	4.07	5.07	6.17	7.30
		98.01	99.97	101.97	104.01	106.09			1.14	1.59	2.20	3.02	4.09
Commodity price			97.03	98.97	100.95	102.97	Commodity		0.28	0.42	0.64	0.97	
				96.06	97.98	99.94	option			0.00	0.00	0.00	
					95.10	97.00					0.00	0.00	
						94.15						0.00	
106.15	107.20	108.26	109.34	110.42	111.51	112.62	4.28	5.21	<b>6.26</b>	<b>7.34</b>	<b>8.42</b>	<b>9.51</b>	10.62
	104.05	105.08	106.12	107.17	108.23	109.30		2.54	3.27	4.15	<b>5.17</b>	<b>6.23</b>	7.30
		101.99	103.00	104.02	105.05	106.09			1.15	1.61	2.22	<b>3.05</b>	4.09
Futures price			99.97	100.96	101.96	102.97	Futures		0.28	0.42	0.64	0.97	
				97.99	98.96	99.94	option			0.00	0.00	0.00	
					96.05	97.00					0.00	0.00	
						94.15						0.00	

**FIGURE 14.10 Lattices associated with a commodity.** The upper left lattice is the price lattice of a commodity. All other lattices are computed from it by backward risk-neutral evaluation.

zero is  $\$100(1.01)^6 = \$106.15$ , as shown in the lattice. The futures price for any node in the lattice can be found by the same technique: just multiply the corresponding commodity price by the factor of interest rate growth for the remaining time.

The futures price can also be found recursively by using the risk-neutral probabilities. We know that the final futures price, at month 6, must be identical to the price of the commodity itself at that time, so we can fill in the last column of the array with those values. Let us denote the futures price at the top of the previous column, at time 5, by  $F$ . If one took the long side of a one-period contract with this assigned price, the payoff in the next period would be either  $112.62 - F$  or  $109.30 - F$ , depending on which of the two nodes was attained. These two values should be multiplied by  $q$  and  $1 - q$ , respectively, and the sum discounted one period to find the initial value, at time 5, of such a contract. But since futures contracts are arranged so that the initial value is zero, it follows that  $q(112.62 - F) + (1 - q)(109.30 - F) = 0$ , which gives  $F = q112.62 + (1 - q)109.30$ . In other words,  $F$  is the weighted average of the next period's prices; the weighting coefficients are the risk-neutral probabilities. We do *not* discount the average.

This process is continued backward a column at a time, computing the weighted average (or expected value) using the risk-neutral probabilities. The final result is again 106.15.

Notice that the original commodity price lattice also can be reconstructed backward by using risk-neutral pricing. Given the final prices, we compute the expected values using the risk-neutral probabilities, but now we *do* discount to find the value at the previous node. Working backward we fill in the entire lattice, duplicating the original figures in the upper left-hand corner.

The backward process for calculating the futures prices and the backward process for computing the commodity prices are identical, except that no discounting is applied in the calculation of futures prices. Hence futures prices will be the same as the commodity prices, but inflated by interest rate growth.

**Example 14.6 (Some options)** Now let us consider some options related to the commodity in Example 14.5. First let us consider a call option on the commodity itself, with a strike price of \$102 and expiration in month 6. This is now easy for us to calculate using binomial lattice methodology, as shown in the upper right-hand part of Figure 14.9. We just fill in the final column and then work backward with the risk-neutral discounting process. The fair price of the option is \$4.16.

Next let us consider a call option on a futures contract with a strike price of \$102. If this option is exercised, the call writer must deliver a futures contract with a futures price of \$102, but marked to market. Suppose the actual futures price at the time of exercise is \$110.42 (as at the top node of the column marked 4 in the futures price lattice). Then the writer can purchase the futures contract (at zero cost) with the futures price \$110.42 and deliver this contract together with the difference of  $\$110.42 - \$102.00 = \$8.42$  to the option holder. This payment compensates for the fact that the writer is delivering a contract at \$110.42, instead of at \$102.00 as promised. In other words, if the option is exercised, the call holder obtains a current futures contract and cash equal to the difference between the current futures price and the option strike price.

We can compute the value of such a call in the same manner as other calls, as shown in the lattice in the lower right-hand portion of Figure 14.9. At each node we must check whether or not it is desirable to exercise the option. This is done by seeing whether the corresponding futures price minus the strike price is greater than the discounted risk-neutral value that would be obtained by holding the option. If it is optimal to exercise the option, we record the option value in boldface. For example, at the top node in the column marked 4, the discounted risk-neutral value is computed to be 8.33. However, by exercising, the option price is found to be \$8.42. Notice that even though the final payoff values are identical for the two options, the futures option has a higher value because the higher intermediate futures prices lead to the possibility of early exercise.

## 14.8 Evaluating Real Investment Opportunities

Options theory can be used to evaluate investment opportunities that are not pure financial instruments. We shall illustrate this by again considering our gold mine lease problems. Now, however, the price of gold is assumed to fluctuate randomly, and this fluctuation must be accounted for in our evaluation of the lease prospect.

**Example 14.7 (Simplico gold mine)** Recall the Simplico gold mine from Chapter 2. Gold can be extracted from this mine at a rate of up to 10,000 ounces per year at a cost of \$200 per ounce. Currently the market price of gold is \$400 per ounce, but we recognize that the price of gold fluctuates randomly. The term structure

0	1	2	3	4	5	6	7	8	9	10
400.0	480.0	576.0	691.2	829.4	995.3	1194.4	1433.3	1719.9	2063.9	2476.7
	360.0	432.0	518.4	622.1	746.5	895.8	1075.0	1289.9	1547.9	1857.5
		324.0	388.8	466.6	559.9	671.8	806.2	967.5	1161.0	1393.1
			291.6	349.9	419.9	503.9	604.7	725.6	870.7	1044.9
				262.4	314.9	377.9	453.5	544.2	653.0	783.6
Gold price (dollars)					236.2	283.4	340.1	408.1	489.8	587.7
						212.6	255.1	306.1	367.3	440.8
							191.3	229.6	275.5	330.6
								172.2	206.6	247.9
									155.0	186.0
										139.5

**FIGURE 14.11 Gold price lattice.** Each year the price either increases by a factor of 1.2 or decreases by a factor of .9. The resulting possible values each year are shown in spreadsheet form.

of interest rates is assumed to be flat at 10%. As a convention, we assume that the price obtained for gold mined in a given year is the price that held at the beginning of the year; but all cash flows occur at the end of the year. We wish to determine the value of a 10-year lease of this mine.

We represent future gold prices by a binomial lattice. Each year the price either increases by a factor of 1.2 (with probability .75) or decreases by a factor of .9 (with probability .25). The resulting lattice is shown in Figure 14.11.

How do we solve the problem of finding the lease value by the methods developed for options pricing? The trick is to notice that the gold mine lease can be regarded as a financial instrument. It has a value that fluctuates in time as the price of gold fluctuates. Indeed, the value of the mine lease at any given time can only be a function of the price of gold and the interest rate (which we assume is fixed). In other words, the lease on the gold mine is a derivative instrument whose underlying security is gold. Therefore the value of the lease can be entered node by node on the gold price lattice.

The lease values on the lattice are determined easily for the final nodes, at the end of the 10 years: the values are zero there because we must return the mine to the owners. At a node representing 1 year to go, the value of the lease is equal to the profit that can be made from the mine that year, discounted back to the beginning of the year. For example, the value at the top node for year 9 is  $10,000(2,063.9 - 200)/1.1 = 16.94$  million. For an earlier node, the value of the lease is the sum of the profit that can be made that year and the risk-neutral expected value of the lease in the next period, both discounted back one period. The risk-neutral probabilities are  $q = (1.1 - .9)/(1.2 - .9) = \frac{2}{3}$ , and  $1 - q = \frac{1}{3}$ . The lease values can therefore be calculated by backward recursion using these values. (At nodes where the price of gold is less than \$200, we do not mine.) The resulting values are indicated in Figure 14.12. We conclude that the value of the lease is \$24,074,548 (showing all the digits).

0	1	2	3	4	5	6	7	8	9	10
24.1	27.8	31.2	34.2	36.5	37.7	37.1	34.1	27.8	16.9	0.0
	17.9	20.7	23.3	25.2	26.4	26.2	24.3	20.0	12.3	0.0
		12.9	15.0	16.7	17.9	18.1	17.0	14.1	8.7	0.0
			8.8	10.4	11.5	12.0	11.5	9.7	6.1	0.0
				5.6	6.7	7.4	7.4	6.4	4.1	0.0
Lease value (millions)					3.2	4.0	4.3	3.9	2.6	0.0
						1.4	2.0	2.1	1.5	0.0
							0.4	0.7	0.7	0.0
								0.0	0.1	0.0
									0.0	0.0
										0.0

**FIGURE 14.12 Simplico gold mine.** The value of the lease is found by working backward. If the price of gold is greater than \$200 per ounce, it is profitable to mine; otherwise no mining is undertaken.

Many readers will be able to see from this example that they have a deeper understanding of investment than they did when they began to study this book. Earlier, in Chapter 2, we discussed the Simplico gold mine under the assumption that the price of gold would remain constant at \$400 over the course of the lease. We also assumed a constant 10% interest rate. These assumptions, which are fairly commonly employed in problems of this type, were probably not regarded as being seriously incongruous by most readers. Now, however, we see that they are not just a simplification, but an actual inconsistency. If the price of gold were known to be constant, gold would act as a risk-free asset with zero rate of return. This is incompatible with the assumption that the risk-free rate is 10%. Indeed, in our lattice of gold prices we must select  $u$ ,  $d$ , and  $R$  such that  $u > R > d$ .

Now that we have “mastered” the Simplico gold mine, it is time to move on to even greater challenges. (If you think you have really mastered the Simplico mine, try Exercise 8.)

**Example 14.8 (Complexico gold mine\*)<sup>3</sup>** The Complexico gold mine was discussed in Chapter 5. In this mine the cost of extraction depends on the amount of gold remaining. Hence if you lease this mine, you must decide how much to mine each period, taking into account that mining in one period affects future mining costs. We also assume now that the price of gold fluctuates according to the binomial lattice of the previous example.

The cost of extraction in any year is \$500  $z^2/x$ , where  $x$  is the amount of gold remaining at the beginning of the year and  $z$  is the amount of gold extracted in ounces. Initially there are  $x_0 = 50,000$  ounces of gold in the mine. We again assume that the term structure of interest rates is flat at 10%. Also, the profit from mining is determined

<sup>3</sup> This is a more difficult example, which should be studied only after you are fairly comfortable with the material of this chapter.

on the basis of the price of gold at the beginning of the period, and in this example all cash flows occur at the beginning of the period.

To solve this problem we must do some preliminary analysis. At the final time the value of the lease is clearly zero. If we are at a node representing the end of year 9, we must determine the optimal amount of gold to mine during the tenth year. Accordingly, we must compute the profit

$$V_9(x_9) = \max_{z_9} (gz_9 - 500z_9^2/x_9)$$

where  $g$  is the price of gold at that particular node. From the calculations of Example 5.5 we know that the maximization gives

$$V_9(x_9) = \frac{g^2 x_9}{2,000}.$$

This shows that the value of the lease is proportional to  $x_9$ , the amount of gold remaining. We therefore write  $V_9(x_9) = K_9 x_9$ , where

$$K_9 = \frac{g^2}{2,000}.$$

We set up a lattice of  $K$  values with nodes corresponding to various gold prices. We put  $K_{10} = 0$  for all elements in the last column and put the values of  $K_9$  in the ninth column. In a similar way, following the analysis of the earlier example, we find that for a node at time 8,

$$V_8(x_8) = \max_{z_8} [gz_8 - 500z_8^2/x_8 + d\hat{K}_9 \times (x_8 - z_8)],$$

where

$$\hat{K}_9 = qK_9 + (1-q)K'_9$$

and where  $K_9$  is the value on the node directly to the right and  $K'_9$  is the value on the node just below that. This leads to

$$z_8 = \frac{(g - d\hat{K}_9)x_8}{1,000}$$

and  $V_8(x_8) = K_8 x_8$ , where

$$K_8 = \frac{(g - \hat{K}_9/R)^2}{2,000} + \hat{K}_9/R.$$

Again, there will be a different value of  $K_8$  for each node at period 8. We work backward with this same formula to complete the lattice shown in Figure 14.13, obtaining  $K_0 = 324.4$ . The value of the lease is then found as  $V_0 = 50,000 \times K_0 = \$16,220,000$ .

0	1	2	3	4	5	6	7	8	9	10
324.4	393.8	478.1	580.8	706.6	862.3	1058.7	1313.4	1656.1	2129.9	0.0
	272.5	329.9	398.6	480.7	578.4	694.4	831.7	995.0	1198.0	0.0
		225.8	272.2	327.0	390.7	463.4	542.9	621.9	673.9	0.0
			182.8	218.9	260.0	305.2	351.1	387.3	379.1	0.0
				143.6	169.5	197.0	222.5	237.3	213.2	0.0
<i>K</i> -value					108.1	124.4	138.1	142.8	119.9	0.0
						76.9	84.1	84.6	67.5	0.0
							50.3	49.5	37.9	0.0
								28.7	21.3	0.0
									12.0	0.0
										0.0

**FIGURE 14.13 Complexico gold mine solution.** The value of the mine is proportional to the amount of gold remaining in the mine. The proportionality factor  $K$  is found by backward recursion.

## Real Options

Sometimes options are associated with investment opportunities that are not financial instruments. For example, when operating a factory, a manager may have the option of hiring additional employees or buying new equipment. As another example, if one acquires a piece of land, one has the option to drill for oil, and then later the option of extracting oil if oil is found. In fact, it is possible to view almost any process that allows control as a process with a series of operational options. These operational options are often termed **real options** to emphasize that they involve *real* activities or *real* commodities, as opposed to purely financial commodities, as in the case, for instance, of stock options. The term *real option* when applied to a general investment problem is also used to imply that options theory can (and should) be used to analyze the problem.

**Example 14.9 (A plant manager's problem)** Some manufacturing plants can be described by a **fixed cost** per month (for equipment, management, and rent) and a **variable cost** (for material, labor, and utilities) that is proportional to the level of production. The total cost is therefore  $T = F + Vx$ , where  $F$  is the fixed cost,  $V$  is the rate of variable cost, and  $x$  is the amount of product produced. The profit of the plant in a month in which it operates at level  $x$  is  $\pi = px - F - Vx$ , where  $p$  is the market price of its product. Clearly, if  $p > V$ , the firm will operate at  $x$  equal to the maximum capacity of the plant; if  $p < V$ , it will not operate. Hence the firm has a continuing option to operate, with a strike price equal to the rate of variable cost. (The Simplico gold mine in Example 14.7 is of this type.)

Real options usually can be analyzed by the same methods used to analyze financial options. Specifically, one sets up an appropriate representation of uncertainty, usually with a binomial lattice, and works backward to find the value. This

solution process is really more fundamental than its particular application to options, so it seems unnecessary and sometimes artificial to force all opportunities for control into options—real or otherwise. Instead, the seasoned analyst takes problems as they come and attacks them directly.

The Simplico mine can be used to illustrate a complex real option associated with the timing of an investment.

**Example 14.10 (Enhancement of the Simplico mine\*)** Recall that the Simplico mine is capable of producing 10,000 ounces of gold per year at a cost of \$200 per ounce. This mine already consists of a whole series of real options—namely, the yearly options to carry out mining operations. In fact, the value of the lease can be expressed as a sum of the values of these individual options (although this viewpoint is not particularly helpful). In this example we wish to consider another option, which is truly in the spirit of a real option.

Suppose that there is a possibility of enhancing the production rate of the Simplico mine by purchasing a new mining machine and making some structural changes in the mine. This enhancement would cost \$4 million but would raise the mine capability by 25% to 12,500 ounces per year, at a total operating cost of \$240 per ounce.

This enhancement alternative is an option, since it need not be carried out. Furthermore, it is an option that is available throughout the term of the lease. The enhancement can be undertaken (that is, exercised) at the beginning of any year, and once in place it applies to all future years. We assume, however, that at the termination of the lease, the enhancement becomes the property of the original mine owner.

Figure 14.14 shows how to calculate the value of the lease when the enhancement option is available. We first calculate the value of the lease assuming that the enhancement is already in place. This calculation is made by constructing the upper lattice of the figure, using exactly the same technique used for the Simplico mine of Example 14.7, but with the new capacity and operating cost figures. The value of the mine under these conditions is \$27.0 million. This figure does not include the cost of the enhancement, so if we were to enhance the mine at time zero, the net value of the lease would be \$23.0 million, which is somewhat less than the value of \$24.1 found earlier without the enhancement. Hence it is not useful to carry out the enhancement immediately.

To find the value of the enhancement option, we construct another lattice, as shown in the lower part of the figure. Here we use the original parameters for production capability and operating cost: 10,000 ounces per year and \$200 per ounce. However, at each node, in addition to the usual calculation of value, we see if it would be useful to jump up to the upper lattice by paying \$4 million. Specifically, we first calculate the value at a node in the lower lattice in the normal way using risk-neutral probabilities. Then we compare this value with the value at the corresponding node in the upper lattice minus \$4 million. We then put the larger of these two values at the node in the lower lattice.

The figures in boldface type show nodes where it is advantageous to jump to the upper lattice by carrying out the enhancement. Note that these values are exactly \$4 million less than their upper counterparts.

**FIGURE 14.14 Option to enhance mine operation.** The top array is computed just as for the Simplico mine, but with parameters of enhancement. The lower array refers to the top one to determine when to carry out the enhancement.

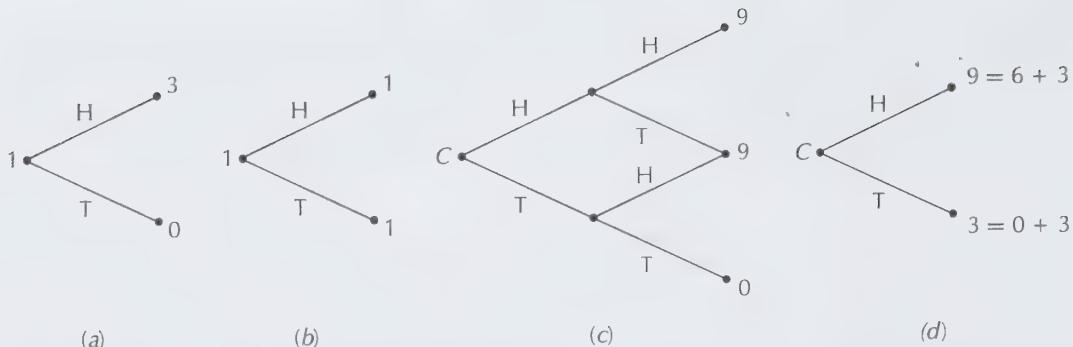
0	1	2	3	4	5	6	7	8	9	10
27.0	31.8	36.4	40.4	43.5	45.2	44.8	41.4	33.9	20.7	0.0
19.5	23.3	26.6	29.3	31.0	31.2	29.2	24.1	14.9	0.0	
	13.5	16.3	18.7	20.4	21.0	20.0	16.8	10.5	0.0	
		8.6	10.8	12.5	13.4	13.2	11.3	7.2	0.0	
			4.9	6.5	7.7	8.0	7.2	4.7	0.0	
Lease value				2.3	3.4	4.1	4.1	2.8	0.0	
assuming enhancement					0.8	1.3	1.8	1.4	0.0	
in place						0.1	0.2	0.4	0.0	
							0.0	0.0	0.0	
								0.0	0.0	
									0.0	
24.6	28.6	32.6	<b>36.4</b>	<b>39.5</b>	<b>41.2</b>	<b>40.8</b>	<b>37.4</b>	<b>29.9</b>	16.9	0
18.0	20.9	23.5		25.6	<b>27.0</b>	<b>27.2</b>	<b>25.2</b>	<b>20.1</b>	12.3	0
12.9	15.0	16.7		17.9	18.1	17.0	14.1	8.7	0	
	8.8	10.4		11.5	12.0	11.5	9.7	6.1	0	
		5.6		6.7	7.4	7.4	6.4	4.1	0	
Lease with option				3.1	4.0	4.3	3.9	2.6	0	
for enhancement					1.3	2.0	2.1	1.5	0	
						0.0	0.7	0.7	0	
						0.0	0.1	0	0.0	
							0.0	0		

The overall value of the lease with the option is given by the value at the first node, and the \$4 million is already taken out. Hence the value of the lease with the enhancement option is \$24.6 million—a slight improvement over the original value of \$24.1 million.

## Linear Pricing

Although we generally use risk-neutral pricing to evaluate derivative securities, it is important to recognize that this evaluation is based on linear pricing; that is, we match a particular derivative to securities we know and then add up the values. The following example highlights the basic simplicity of the method.

**Example 14.11 (Gavin explains)** Gavin was excited by the new pricing method, and he was anxious to show his father, D. Jones, how it worked. He took a quarter out his pocket and explained the setup. “You can bet on the outcome of a coin flip. If you bet a dollar and heads comes up, you get \$3, but if the outcome is tails, you get nothing and lose your dollar. This is like a risky stock. Of course you can keep the coin in your pocket; then after the flip you will still have your dollar. Your pocket is a risk-free alternative that we may call a bond.” When his father nodded that he understood, Gavin continued. “You can bet at any level, even negative to short the coin flip. And you can do either at any time.”



**FIGURE 14.15 A proposition and its parts.** Tree (a) is a basic risky proposition; tree (b) is a risk-free opportunity; and tree (c) represents a new, more complex proposition. The value  $C$  can be found by breaking it into its parts. The final piece is shown in (d).

Gavin flipped the coin a few times just to highlight the situation. He continued, “Now, here is a new proposition. I flip the coin twice, and if at least one outcome is heads, you will get \$9. If both flips are tails, you get nothing. Of course you can also bet on each flip individually as before if you wish. How much is this new proposition worth?”

His father thought a moment and then said, “I could work out the probabilities.”  
Gavin said, “It has nothing to do with actual probabilities.”

Mr. Jones looked doubtful. “I guess you better show me.”

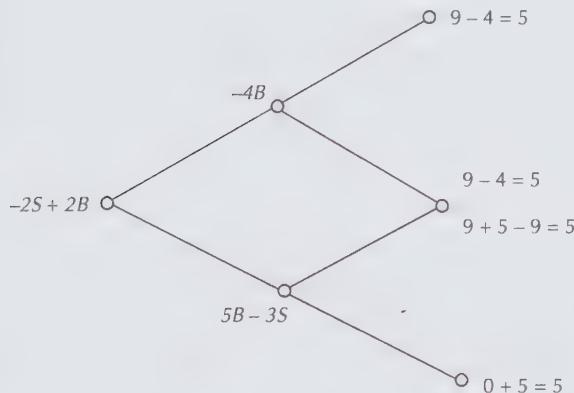
On a piece of paper, Gavin drew the first three trees of Figure 14.15, showing the original possibilities and the new proposition. Then he asked how much it was worth to be at the position where a single flip had occurred and it was a head. Mr. Jones quickly understood that it was worth \$9 because the next step was like 9 times the risk-free payoff. Likewise, he figured that if the first flip was a tail, the payoff from there was the same as 3 times the stock payoff. With Gavin’s help he saw that the payoff of 9 and 3 could be broken into 6 and 0 plus 3 and 3, for a total of  $2 + 3 = 5$ .

Gavin then showed that using expected value would give the correct answer provided one used the risk-neutral probability of  $1/3$ .

**Example 14.12 (Perfect hedging)** The derivation of the price of 5 for the coin flip example should be convincing, but there is still the issue of risk. If someone were to promise you the payoff of the two-flip example, you would apparently have risk, since you do not actually get 5 for certain, but, rather, you get either 9 or 0.

However, it is in fact possible to transform that uncertain payoff to one that gives 5 for sure and maintains a zero balance at every intermediate step. The trick is to hedge appropriately—and the result of this hedge should further convince one that the true price is 5.

Consider the lattice of Figure 14.16. You begin at the initial point by shorting two units of the stock (the up move) and go long two units of the risk-free asset. The



**FIGURE 14.16** A proper hedge will produce a risk-free result.

cost is zero. In the figure, the first node is labeled  $-2S + 2B$  to represent this action. Then if a head occurs, you will have a payoff of  $-6 + 2 = -4$ . That means that you owe \$4. To clear that current balance, you borrow 4 units of the bond, as indicated. If at the next flip you get either a head or a tail, you will get \$9 from the proposition and pay back your loan of 4 to end up with 5. The upper calculation is used at the middle node.

Likewise, if the first flip is a tail, you will have \$2 at the resulting node. You will then go long 5 units of the bond and short 2 units of the stock. Including the 2 you already have, this strategy gives you a 0 current balance. Finally, if the next flip is a head, you end up with  $9 - 9 + 5 = 5$ , from the lower calculation at the middle node, whereas if the outcome is a tail, you will end up with  $0 + 5 = 5$ .

No matter what happens you get 5!

## 14.9 General Risk-Neutral Pricing\*

A general principle of risk-neutral pricing can be inferred from the analysis and methods of the previous few sections. This principle provides a compact formula for the price of a derivative security under the binomial lattice formulation.

Suppose that the price  $S$  of an asset is described by a binomial lattice, and suppose that  $f$  is a security whose cash flow at any time  $k$  is a function only of the node at time  $k$ . Then the arbitrage-free price of the asset is

$$f_{va1} = \hat{E} \left( \sum_{k=0}^N d_k f_k \right). \quad (14.12)$$

In this equation the summation represents the discounted cash flow, with the  $d_k$ 's being the risk-free discount factors as seen at time 0. The  $f_k$ 's are the period cash flows, which depend on the particular node at  $k$  that occurs. Hence the  $f_k$ 's are random. The

expectation  $\hat{E}$  is taken with respect to the risk-neutral probabilities associated with the lattice of the underlying asset.

Consider a European call option with strike price  $K$ . The pricing formula, Eq. (14.12), becomes

$$C = \frac{1}{R_T} \hat{E}[\max(S_T - K, 0)], \quad (14.13)$$

where  $R_T$  is the risk-free return for the whole time to expiration. In this case there is only a single cash flow,  $\max(S_T - K, 0)$ , occurring at the final time. We take the risk-neutral expected value of this and discount it to the present. Note that actual calculation using this formula is best done by working backward from the end. We use the running present value method to back the formula up one stage at a time.

In many situations the cash flow stream can be influenced by our actions as well as by chance. For instance, we may have the opportunity to exercise an option before expiration, decide how much gold to mine, or add enhancements. In such cases the general pricing formula becomes

$$f_{\text{val}} = \max \left[ \hat{E} \left( \sum_{k=0}^N d_k f_k \right) \right],$$

where the maximization is taken with respect to the available actions. We have seen in the examples of this chapter how this maximization can in many cases be carried out as part of the backward recursion process, although the size of the lattice sometimes must be increased. This general formula has great power, for it provides a way to formulate and solve many interesting and important investment problems.

## 14.10 Three-principle Power

Consider again the Simplico Gold Mine of Example 14.7. We suggest that it is possible to find its value without constructing even a single lattice, but instead making only a couple of simple calculations based on the three pricing principles. You might try to see how to do it before proceeding.

Each year there is an option to mine or not, but it is always optimal to mine unless the price of gold falls below \$200/ounce, which only occurs in a few cells at the bottom right of the price lattice. Ignoring this option by mining each year lowers the value by only about \$25 (depending on numerical accuracy).

The mine income comes from the sale of 10,000 ounces of gold each period. These prices are unknown at time 0, but from the first pricing principle the value at time 0 of the sale of an ounce of gold at a later period is simply the current price, \$400.00. Hence, the total value of all receipts of gold sales is  $10,000 \times 10 \times \$400 = \$40,000,000$ . Actually, since the money is obtained at the end of each period, we must discount by one period of 10% interest. Thus the total value, at time 0, of all gold receipts is  $G = \$400,000,000 / 1.1 = 36,363,636.36$ .

The costs are \$200 per ounce each period, and these are fixed. Hence, from the second pricing principle, the present value of the total cost is the value of a 10-year

annuity of \$200 each year, times 10,000. We can find this easily on a spreadsheet or by use of the annuity formula from Section 3.2, which is  $C = 10,000 \frac{\$200}{.10} \left[ 1 - \frac{1}{1.1^{10}} \right] = \$12,289,134.21$ . The final value of the mine at time 0 is  $G - C = 24,074,502.15$ , which agrees almost exactly with the lattice method. The linearity principle was used to combine the various individual cash flows.

You might try Exercise 8 now.

## Decomposition of the Pricing Principles

The fundamental pricing principles are directly applicable when the payoff of an asset is a linear function of an underlying marketed asset. This is the case for some derivatives such as forwards, but for others, such as options, it is not. However, the principles do apply to an asset payoff determined by a single binomial step, since any outcome can be expressed as a linear combination of the underlying stock and bond. A multistep binomial process (a lattice or tree), can be converted to a series of linear steps, so that the fundamental principles can price the entire process. In the next chapter we shall see how the fundamental principles can in a similar way be applied in continuous time as well.

## 14.11 Summary

An option is the right, but not the obligation, to buy (or sell) an asset under specified terms. Options have had a checkered past, but for the past three decades they have played an important role in finance. Used wisely, they can control risk and enhance the performance of a portfolio. Used carelessly, options can greatly increase risk and lead to substantial losses.

Options terminology includes: call, put, exercise, strike price, expiration, writing a call, premium, in the money, out of the money, American option, and European option.

A major topic of options theory is the determination of the correct price (or premium) of an option. This price depends on the price of the underlying asset, the strike price, the time to expiration, the volatility of the underlying asset, the cash flow generated by the asset (such as dividend payments), and the prevailing interest rate. Although determination of an appropriate option price can be difficult, certain relations can be derived from simple no-arbitrage arguments. For example, for European-style options there is parity between a put and a call with the same strike price. Likewise, the value of a combination of options (such as in a butterfly spread) must be the same combination of the prices of the component options.

One important result is that it is never optimal to exercise, before expiration, an American call option on a stock that does not pay a dividend before expiration.

A general way to find the price of an option is to use the binomial lattice methodology. The random process of the underlying asset is modeled as a binomial lattice. The value of the option at expiration is entered on the final nodes of a corresponding option lattice. The other nodes in the option lattice are computed one at a time

by working backward through the periods. For a European-style option (without the possibility of early exercise) the value at any node in the option lattice is found by computing the expected value of the value next period using risk-neutral probabilities. This expected value is then discounted by the effect of one period's interest rate. If the option is an American-style option, the value computed as before must be compared with the value that could be obtained by exercise at that time, and the greater of the two compared values is taken to be the final value for that node.

The risk-neutral probabilities are easy to calculate. The risk-neutral probability for an up move is  $q = (R - d)/(u - d)$ . The easiest way to derive this formula is to find the  $q$  that makes the price of the underlying security equal to the discounted expected value of its next period value.

The binomial lattice methodology can be used to find the value of other investments besides options. Indeed, it can be used to evaluate any project whose cash flow stream is determined by an underlying traded asset. Examples include futures on options, gold mine leases, oil wells, and tree farms. With ingenuity, even complex real options can be evaluated by constructing two or more interrelated binomial lattices.

## Exercises

- (Bull spread)** An investor who is bullish about a stock (believing that it will rise) may wish to construct a *bull spread* for that stock. One way to construct such a spread is to buy a call with strike price  $K_1$  and sell a call with the same expiration date but with a strike price of  $K_2 > K_1$ . Draw the payoff curve for such a spread. Is the initial cost of the spread positive or negative?
- (Put-call parity)** Suppose over the period  $[0, T]$  a certain stock pays a dividend whose present value at interest rate  $r$  is  $D$ . Show that the put-call parity relation for European options at  $t = 0$ , expiring at  $T$ , is

$$C + D + Kd = P + S,$$

where  $d$  is the discount factor from 0 to  $T$ .

- (Patent insurance)** One year ago, Bioette, a biotech incubator, entered into a forward contract to sell one of its patents to Pharm, a major drug company, in 2 years for \$10 million.

Currently, with only 1 year remaining in the contract, both parties are concerned about the uncertainty of the patent's market value. A major insurer, Protech, is offering loss insurance to both parties. For a premium of \$.4 million now, it will pay Bioette the difference between the market value of the patent and the fixed sales price of \$10 million if the market price is larger. Otherwise, it will pay nothing.

In a similar way, for a premium of \$1 million, Protech will pay Pharm the difference between the \$10 million and the market price if the market price is less than \$10 million. Otherwise, it will pay nothing.

Considering that the price of a zero-coupon bond with a face value of \$100 and 1 year to maturity is \$90, what is the value that Protech would assign to the patent? (You may assume the insurance is tradable.)

- (Call strikes  $\diamond$ )** Consider a family of call options on a non-dividend-paying stock, each option being identical except for its strike price. The value of the call with strike price  $K$  is denoted by  $C(K)$ . Prove the following three general relations using arbitrage arguments:

- (a)  $K_2 > K_1$  implies  $C(K_1) \geq C(K_2)$ .
- (b)  $K_2 > K_1$  implies  $K_2 - K_1 \geq C(K_1) - C(K_2)$ .
- (c)  $K_3 > K_2 > K_1$  implies

$$C(K_2) \leq \left( \frac{K_3 - K_2}{K_3 - K_1} \right) C(K_1) + \left( \frac{K_2 - K_1}{K_3 - K_1} \right) C(K_3).$$

- 5. (Fixed dividend  $\oplus$ )** Suppose that a stock will pay a dividend of amount  $D$  at time  $\tau$ . We wish to determine the price of a European call option on this stock using the lattice method. Accordingly, the time interval  $[0, T]$  covering the life of the option is divided into  $N$  intervals, and hence  $N + 1$  time periods are assigned. Assume that the dividend date  $\tau$  occurs somewhere between periods  $k$  and  $k + 1$ . One approach to the problem would be to establish a lattice of stock prices in the usual way, but subtract  $D$  from the nodes at period  $k$ . This produces a tree with nodes that do not recombine, as shown in Figure 14.17.

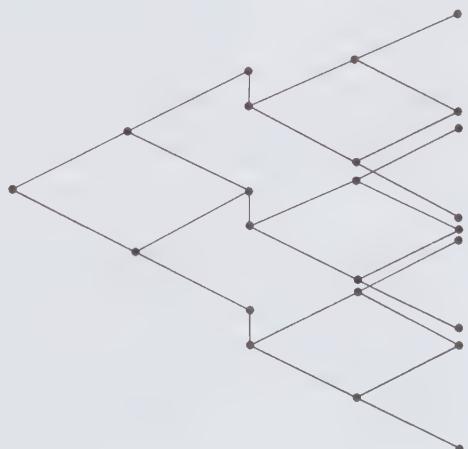
The problem can be solved this way, but there is another representation that does recombine. Since the dividend amount is known, we regard it as a nonrandom component of the stock price. At any time before the dividend we regard the price as having two components: a random component  $S^*$  and a deterministic component equal to the present value of the future dividend. The random component  $S^*$  is described by a lattice with initial value  $S(0) - De^{-r\tau}$  and with  $u$  and  $d$  determined by the volatility  $\sigma$  of the stock. The option is evaluated on this lattice. The only modification that must be made in the computation is that when valuing the option at a node, the stock price used in the valuation formula is not just  $S^*$  at that node, but rather  $S = S^* + De^{-r(\tau-t)}$  for  $t < \tau$ . Use this technique to find the value of a 6-month call option with  $S(0) = 50$ ,  $K = 50$ ,  $\sigma = 20\%$ ,  $R = 10\%$ , and  $D = \$3$  to be paid in  $3\frac{1}{2}$  months.

- 6. (Call inequality)** Consider a European call option on a non-dividend-paying stock. The strike price is  $K$ , the time to expiration is  $T$ , and the price of one unit of a zero-coupon bond maturing at  $T$  is  $B(T)$ . Denote the price of the call by  $C(S, t)$ . Show that

$$C(S, t) \geq \max[0, S - KB(T)].$$

[Hint: Consider two portfolios: (a) purchase one call, (b) purchase one share of stock and sell  $K$  bonds.]

**FIGURE 14.17 Nonrecombining dividend tree.**



7. (Perpetual call) A perpetual option is one that never expires. (Such an option must be of American style.) Use Exercise 6 to show that the value of a perpetual call on a non-dividend-paying stock is  $C = S$ .
8. (A surprise  $\oplus$ ) Consider a deterministic cash flow stream  $(x_0, x_1, x_2, \dots, x_n)$  with all positive flows. Let  $PV(r)$  denote the present value of this stream at an interest rate  $r$ .
- If  $r$  decreases, does  $PV(r)$  increase or decrease?
  - Solve the Simplico gold mine problem with  $r = 4\%$  and find that the value of the lease is \$22.1 million. Can you explain why the value decreased relative to its value with  $r = 10\%$ ?
9. (My coin) There are two propositions: (a) I flip a coin. If it is heads, you are paid \$3; if it is tails, you are paid \$0. It costs you \$1 to participate in this proposition. You may do so at any level, or repeatedly, and the payoffs scale accordingly. (b) You may keep your money in your pocket (earning no interest). Here is a third proposition: (c) I flip the coin three times. If at least two of the flips are heads, you are paid \$27; otherwise zero. How much is this proposition worth?
10. (The happy call) A New York firm is offering a new financial instrument called a “happy call.” It has a payoff function at time  $T$  equal to  $\max(.5S, S - K)$ , where  $S$  is the price of a stock and  $K$  is a fixed strike price. You always get something with a happy call. Let  $P$  be the price of the stock at time  $t = 0$  and let  $C_1$  and  $C_2$  be the prices of ordinary calls with strike prices  $K$  and  $2K$ , respectively. The fair price of the happy call is of the form

$$C_H = \alpha P + \beta C_1 + \gamma C_2.$$

Find the constants  $\alpha$ ,  $\beta$ , and  $\gamma$ .

11. (You are a president) It is August 6. You are the president of a small electronics company. The company has some cash reserves that will not be needed for about 3 months, but interest rates are very low. Your chief financial officer (CFO) tells you that a progressive securities firm has an investment that guarantees no losses and allows participation in upward movements of the stock market. In fact, the total rate of return until the third week of November is to be determined by the formula  $\max(0, .25r)$ , where  $r$  is the rate of return on the S&P 100 stock index during the 3-month period (ignoring dividends). The CFO suggests that this conservative investment might be an ideal alternative to participation in the interest rate market and asks for your opinion. You pick up *The Wall Street Journal* and make a few simple calculations to check whether it is, in fact, a good deal. Show these calculations and the conclusion. Data: A stock index is 14.2. There are two call options for November with strikes 410 and 420 and prices 13 and 7.5. The yield on November Treasury bills is 3.11.
12. (Simplico invariance) If the Simplico mine is solved with all parameters remaining the same except that  $u = 1.2$  is changed to  $u = 1.3$ , the value of the lease remains unchanged to within three decimal places. Indeed, quite wide variations in  $u$  and  $d$  have almost no influence on the lease price. Give an intuitive explanation for this.
13. (Change of period length  $\oplus$ ) A stock has volatility  $\sigma = .30$  and a current value of \$36. An American put option on this stock has a strike price of \$40, and expiration is in 5 months. The interest rate is 8%. Find the value of this put using a binomial lattice with 1-month intervals. Repeat using a lattice with half-month intervals.
14. (Average value Complexico  $\oplus$ ) Suppose that the price received for gold extracted from time  $k$  to  $k+1$  is the average of the price of gold at these two times; that is,  $(g_k + g_{k+1})/2$ .

However, costs are incurred at the beginning of the period whereas revenues are received at the end of the period. Find the value of the Complexico mine in this case.

- 15.** (“As you like it” option) Consider the stock of Examples 14.3 and 14.4, which has  $\sigma = .20$  and an initial price of \$62. The interest rate is 10%, compounded monthly. Consider a 5-month option with a strike price of \$60. This option can be declared, after exactly 3 months, by the purchaser to be either a European call or a European put. Find the value of this “as you like it” option.
- 16.** (Tree harvesting  $\oplus$ ) You are considering an investment in a tree farm. Trees grow each year by the following factors:

Year	1	2	3	4	5	6	7	8	9	10
Growth	1.6	1.5	1.4	1.3	1.2	1.15	1.1	1.05	1.02	1.01

The price of lumber follows a binomial lattice with  $u = 1.20$  and  $d = .9$ . The interest rate is constant at 10%. It costs \$2 million each year, payable at the beginning of the year, to lease the forest land. The initial value of the trees is \$5 million (assuming they were harvested immediately). You can cut the trees at the end of any year and then not pay rent after that. (For those readers who care, we assume that cut lumber can be stored at no cost.)

- (a) Argue that if the rent were zero, you would never cut the trees as long as they were growing.  
 (b) With rent of \$2 million per year, find the best cutting policy and the value of the investment opportunity.
- 17.** (Coin market) There is a market for bets on the outcome of a coin toss. The possible outcomes are heads, tails, and edge. There are three assets traded in that market:

Asset A pays \$1 independent of the outcome.

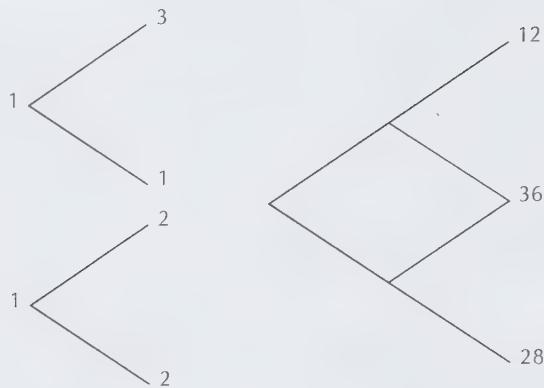
Asset B pays \$1.50 for a head and \$0 for tails and \$1 if the coin lands on its edge.

Asset C pays \$10 if and only if the coin lands on its edge (and \$0 otherwise).

The prices of those assets are always constant and are fixed at \$1. We assume that the payments to the winning bidders take place immediately after the outcome is determined.

- (a) What is the implied risk-free interest rate in this market?  
 (b) What are the risk-neutral probabilities of each of the possible coin toss outcomes (heads, tails, edge)?
- 18.** (High interest) Consider the stock and bond lattices on the left of Figure 14.18. Each has unit cost of \$1 with the payoffs shown at the end nodes. There is a new asset available—a derivative of the first two—with the payoffs shown at the end of two periods, as indicated in the lattice on the right side of the figure.

- (a) Find the value of this new asset.  
 (b) What is the deterministic payoff that should be obtained at the end of two periods?  
 (c) Show how to hedge this new asset by trading the stock and the bond in such a way that the deterministic payoff is obtained. [Hint: Construct a lattice with payoffs 16, -8, and 0, and argue that it should cost 0. Solve for the resulting values at each node. Next work forward using a self-financing portfolio to match the node values. You may like to remember that to match in a binary lattice we set  $x = (Cu - Cd)/(u - d)$ .]

**FIGURE 14.18** Hedge out the uncertainty.

- 19.** (Combination option) A derivative security of European style with expiration in 1 year has this payoff:  $\max(0, \min(3K - S, S - K))$ , where  $K = 10$  is the strike price and  $S$  is the price of the underlying stock at expiration. The stock currently trades at 25, and the following prices for European options on the stock are known (all expiring in 1 year):

Type	Strike	Price
Call	10	16.76
Call	20	7.02
Put	10	0.85
Put	30	6.05

- (a) Draw the graph of the payoff as a function of  $S$ .
- (b) What is the 1-year interest rate  $r$ ?
- (c) What is the price  $P$  of the derivative security?

## References

For general background material on options, see [1–3]. The pricing of options was originally addressed mathematically by Bachelier [4] using a statistical approach. The analysis of put-call parity and various price inequalities that hold independently of the underlying stock process was systematically developed in [5]. The rational option price based on the no-arbitrage principle was first discovered by Black and Scholes [6] when the price of the underlying asset was governed by geometric Brownian motion. The simplified approach using a binomial lattice was first presented in [7] and later developed in [8, 9]. The risk-neutral formulation of option evaluation was generalized to other derivatives in [10]. Exercise 4 is adopted from [2].

1. Sincere, M., (2007), *Understanding Options*, McGraw-Hill, New York.
2. Cox, J. C., and M. Rubinstein (1985), *Options Markets*, Prentice Hall, Englewood Cliffs, NJ.
3. Hull, J. C. (2008), *Options, Futures, and Other Derivative Securities*, 7th ed., Prentice Hall, Englewood Cliffs, NJ.

4. Bachelier, L. (1900), "Théorie de la Spéculation," *Annals de l'Ecole Normale Supérieure*, **17**, 21–86. English translation by A. J. Boness (1967) in *The Random Character of Stock Market Prices*, P. H. Cootner, Ed., M.I.T. Press, Cambridge, MA, 17–78.
5. Merton, R. C. (1973), "Theory of Rational Option Pricing," *Bell Journal of Economics and Management Science*, **4**, 141–183.
6. Black, F., and M. Scholes (1973), "The Pricing of Options and Corporate Liabilities," *Journal of Political Economy*, **81**, 637–654.
7. Sharpe, W. F. (1978), *Investments*, Prentice Hall, Englewood Cliffs, NJ.
8. Cox, J. C., S. A. Ross, and M. Rubinstein (1979), "Option Pricing: A Simplified Approach," *Journal of Financial Economics*, **7**, 229–263.
9. Rendleman, R. J., Jr., and B. J. Bartter (1979), "Two-State Option Pricing," *Journal of Finance*, **34**, 1093–1110.
10. Harrison, J. M., and D. M. Kreps (1979), "Martingales and Arbitrage in Multiperiod Securities Markets," *Journal of Economic Theory*, **20**, 381–408.

# 15

## ADDITIONAL OPTIONS TOPICS

### 15.1 Introduction

Options theory plays a major role in the modern theory of finance because it so clearly highlights the power of the comparison principle, based on the assumption that there are no arbitrage opportunities. The previous chapter presented the theory in a simple and practical form, using the binomial lattice framework. That material is by itself sufficient to solve most options problems. There is, however, a continuous-time version of the theory and extensions of the lattice theory, which lead to new financial insights, allow consideration of more complex derivative securities, provide alternative computational methods, and prepare the way for the more complete theory of investment presented in the following chapters.

### 15.2 The Black–Scholes Equation

The famous Black–Scholes option pricing equation initiated the modern theory of finance based on the no-arbitrage principle. Its development triggered an enormous amount of research and revolutionized the practice of finance. The equation was developed under the assumption that the price fluctuations of the underlying security can be described by an Ito process, as presented in Chapter 13. The logic behind the equation is, however, conceptually identical to that used for the binomial lattice: at each moment two available securities are combined to construct a portfolio that reproduces the local behavior of the derivative security. Historically, the Black–Scholes

theory of options predated the binomial lattice theory by several years, the lattice theory being a result of simplification.

To begin the presentation of the Black–Scholes equation, let the price  $S$  of an underlying security (which we shall refer to as a stock) be governed by a geometric Brownian motion process over a time interval  $[0, T]$  described by

$$dS = \mu S dt + \sigma S dz, \quad (15.1)$$

where  $z$  is standard Brownian motion (or a Wiener process). Suppose also that there is a risk-free asset (a bond) carrying an interest rate of  $r$  over  $[0, T]$ . The value  $B$  of this bond satisfies

$$dB = rB dt. \quad (15.2)$$

Finally consider a security that is derivative to  $S$ , which means that its price is a function of  $S$  and  $t$ . Let  $f(S, t)$  denote the price of this security at time  $t$  when the stock price is  $S$ . We want a (nonrandom) equation for the function  $f(S, t)$ , which will give the price of the derivative explicitly. This function can be found by solving the Black–Scholes equation as stated:

**Black–Scholes equation** *Suppose that the price of a security is governed by the process (15.1) and the interest rate is  $r$ . A derivative of this security has a price  $f(S, t)$ , which satisfies the partial differential equation*

$$\frac{\partial f}{\partial t} + \frac{\partial f}{\partial S} rS + \frac{1}{2} \frac{\partial^2 f}{\partial S^2} \sigma^2 S^2 = rf. \quad (15.3)$$

We present a proof of this result later in this section, but first let us look at its significance.

As a simple example, consider the stock itself. It is (in a trivial way) a derivative of  $S$ , so  $f(S, t) = S$  should satisfy the Black–Scholes equation. In fact, with this choice of  $f$  we have  $\partial f / \partial t = 0$ ,  $\partial f / \partial S = 1$ ,  $\partial^2 f / \partial S^2 = 0$ . Hence equation (15.3) becomes  $rS = rS$ , which shows that  $f(S, t) = S$  is a solution.

As another simple example, consider the bond. It also is (in a trivial way) a derivative of  $S$ , so  $f(S, t) = e^{rt}$  should satisfy the Black–Scholes equation. In fact, with this choice of  $f$  we have  $\partial f / \partial t = re^{rt}$ ,  $\partial f / \partial S = 0$ ,  $\partial^2 f / \partial S^2 = 0$ . Hence equation (15.3) becomes  $re^{rt} = re^{rt}$ , which shows that, indeed,  $f(S, t) = e^{rt}$  is a solution. There are uncountably more solutions.

In general, the Black–Scholes equation can be thought of in two ways. First, suppose that we arbitrarily specify a function  $f(S, t)$  and announce that this is the price of a new security. Since we specify the function, we can arrange for it *not* to satisfy the Black–Scholes equation. What is wrong? If  $f(S, t)$  does not satisfy the Black–Scholes equation, then there is an arbitrage opportunity lying somewhere among  $S$ ,  $B$ , and  $f$ . By a proper combination of these (and the combination may change with time) it will be possible to extract money, risk free. Hence the first way to look at the Black–Scholes equation is that it establishes a property that must hold for a derivative security's price function.

The second way to view the equation is that it can be used to actually find the price function corresponding to various derivative securities. This is done by

specifying appropriate boundary conditions that are used in conjunction with the Black–Scholes partial differential equation to solve for the price function. For example, specifying  $f(S, T) = S(T)$  leads to  $f(S, t) = S(t)$ ; specifying  $f(S, T) = e^{rT}$  leads to  $f(S, t) = e^{rt}$ . As a nontrivial example, the price  $C(S, t)$  of a European call option on a stock that pays no dividends must satisfy the Black–Scholes equation (with  $C$  playing the role of  $f$ ) and it must satisfy the boundary conditions

$$C(0, t) = 0 \quad (15.4)$$

$$C(S, T) = \max(S - K, 0). \quad (15.5)$$

Likewise, for a European put with price  $P(S, t)$  the boundary conditions are

$$P(\infty, t) = 0 \quad (15.6)$$

$$P(S, T) = \max(K - S, 0). \quad (15.7)$$

Other derivative securities may have different forms of boundary conditions, which are sufficient to determine the entire function  $f(S, T)$ . For example, the boundary conditions for an American call option and an American put on a non-dividend-paying stock require, in addition to the conditions mentioned, a condition concerning the possibility of early exercise. These are

$$C(S, t) \geq \max(0, S - K) \quad (15.8)$$

$$P(S, t) \geq \max(0, K - S). \quad (15.9)$$

Of course, the additional boundary condition for calls is unnecessary, since an American call on a non-dividend-paying stock is never exercised early.

**Example 15.1 (A perpetual call)** Consider a perpetual call option with strike price  $K$ . There is no terminal boundary condition since  $T = \infty$ . However, it must be American style, and the no-early-exercise condition  $f(S, t) \geq \max(0, S - K)$  for all  $t$  must be satisfied by the solution  $f$ . In addition, we must have  $f(S, t) \leq S$  for all  $t$  since the call must cost less than the security itself. As an (informed) guess we might try the simple solution  $f = S$ . Indeed, we know that this satisfies the Black–Scholes equation. The two boundary conditions are also satisfied.

The solution  $f(S) = S$  for the value of a perpetual call does make intuitive sense. If the call is held for a long time, the stock value will almost certainly increase to a very large value, so that the exercise price  $K$  is insignificant in comparison. Hence if we owned the call we could obtain the stock later for essentially nothing, duplicating the position we would have if we initially bought the stock.

## Proof of the Black–Scholes Equation\*

How can we derive the Black–Scholes equation? The key idea is the same idea used in Chapter 14 to derive the binomial lattice pricing method. At any time we form a portfolio with portions of the stock and the bond so that this portfolio exactly matches the (instantaneous) return characteristics of the derivative security. The value

of this portfolio must equal the value of the derivative security. In a binomial lattice framework the matching is done period by period, relating the value at one time point to those at the next. In the continuous-time framework, the matching is done at each instant, relating the value at one time to the rates of change at that time. Replication is used in both cases. Here is the conventional proof. (See Section 15.13 for an alternative.)

**Proof:** By Ito's lemma [equation (13.22)] we have

$$df = \left( \frac{\partial f}{\partial t} + \frac{\partial f}{\partial S} \mu S + \frac{1}{2} \frac{\partial^2 f}{\partial S^2} \sigma^2 S^2 \right) dt + \frac{\partial f}{\partial S} \sigma S dz, \quad (15.10)$$

which is an Ito process for the price of the derivative security. This price fluctuates randomly along with the stock price  $S$  and the Brownian motion  $z$ .

We form a portfolio of  $S$  and  $B$  that replicates the behavior of the derivative security. In particular, at each time  $t$  we select an amount  $x_t$  of the stock and an amount  $y_t$  of the bond, giving a total portfolio value of  $G(t) = x_t S(t) + y_t B(t)$ . We wish to select  $x_t$  and  $y_t$  so that  $G(t)$  replicates the derivative security value  $f(S, t)$ . The instantaneous gain in value of this portfolio due to changes in security prices (the investment gain) is

$$dG = x_t dS + y_t dB. \quad (15.11)$$

Expanding, we write

$$\begin{aligned} dG &= x_t dS + y_t dB \\ &= x_t \{ \mu S dt + \sigma S dz \} + y_t rB dt \\ &= (x_t \mu S + y_t rB) dt + x_t \sigma S dz. \end{aligned} \quad (15.12)$$

Since we want the portfolio gain of  $G(t)$  to behave just like the gain of  $f$ , we match the coefficients of  $dt$  and  $dz$  in equation 15.12 to those of (15.10). To do this we first match the  $dz$  coefficient by setting

$$x_t = \frac{\partial f}{\partial S}. \quad (15.13)$$

Requiring  $G = x_t S + y_t B$  and  $G = f$  gives

$$y_t = \frac{1}{B} \left[ f(S, t) - S \frac{\partial f}{\partial S} \right].$$

Substituting these expressions in equation (15.12) and matching the coefficient of  $dt$  in equation (15.10), we obtain

$$\frac{\partial f}{\partial S} \mu S + \frac{1}{B} \left[ f(S, t) - S \frac{\partial f}{\partial S} \right] rB = \frac{\partial f}{\partial t} + \frac{\partial f}{\partial S} \mu S + \frac{1}{2} \frac{\partial^2 f}{\partial S^2} \sigma^2 S^2.$$

Or, finally,

$$\frac{\partial f}{\partial t} + \frac{\partial f}{\partial S} rS + \frac{1}{2} \frac{\partial^2 f}{\partial S^2} \sigma^2 S^2 = rf. \quad (15.14)$$

This is the Black-Scholes equation. ■

## Self-Financing Strategies\*

A portfolio strategy controls the composition of a portfolio, and there are two basic ways that the portfolio value  $G$  can change: by changes in the market values of  $S$  and  $B$ , and through infusion or withdrawal of capital. The change due to purely market changes is  $dG(t) = x_t dS(t) + y_t dB(t)$ . Basically, for given levels of  $x_t$  and  $y_t$ , the portfolio value changes due to changes in the prices of the market assets. Alternatively, value can change if market values stay fixed but the investment levels  $x_t$  and  $y_t$  are changed in such a way that their total expense changes. However, in the proof of the Black–Scholes equation, the construction of the replicating portfolio imposes the constraint (15.11). That constraint guarantees that no infusion or withdrawal of cash is used. We say that the replicating strategy is **self-financing**, since the value  $G(S, t)$  of the replicating portfolio matches  $f(S, t)$  with changes in market value only.<sup>1</sup>

## 15.3 Call Option Formula

Although it is usually impossible to find an analytic solution to the Black–Scholes equation, it is possible to find such a solution for a European call option. This analytic solution is of great practical and theoretical use.

The formula uses the function  $N(x)$ , the standard **cumulative normal probability distribution**. This is the cumulative distribution of a normal random variable having mean 0 and variance 1. It can be expressed as

$$N(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-y^2/2} dy. \quad (15.15)$$

The function  $N(x)$  is illustrated in Figure 15.1. The value  $N(x)$  is the area under the familiar bell-shaped curve from  $-\infty$  to  $x$ . Particular values are  $N(-\infty) = 0$ ,  $N(0) = \frac{1}{2}$ , and  $N(\infty) = 1$ .

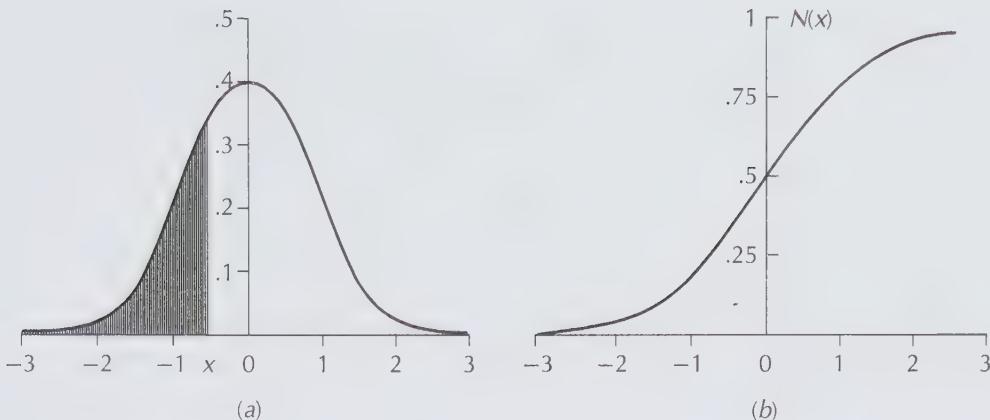
The function  $N(x)$  cannot be expressed in closed form, but there are tables for its values, and there are accurate approximation formulas. (See Exercise 1.)

**Black–Scholes call option formula** Consider a European call option with strike price  $K$  and expiration time  $T$ . If the underlying stock pays no dividends during the time  $[0, T]$  and if interest is constant and continuously compounded at a rate  $r$ , the Black–Scholes solution is  $f(S, t) = C(S, t)$ , defined by

$$C(S, t) = SN(d_1) - Ke^{-r(T-t)}N(d_2), \quad (15.16a)$$

---

<sup>1</sup> In general, the total change in  $G$  is  $dG(S, t) = x_t dS(t) + y_t dB(t) + dx_t S(t) + dy_t B(t) + dx_t dS(t) + dy_t dB(t)$ , but it can be shown that the sum of the last four terms is zero under the Black–Scholes construction.



**FIGURE 15.1 Normal density and cumulative distribution.** (a) The curve is the normal density  $(1/\sqrt{2\pi})e^{-x^2/2}$ . The area under the curve up to the point  $x$  gives the value of the cumulative distribution  $N(x)$ . (b) The cumulative distribution itself rises smoothly from 0 to 1, but it does not have a closed-form representation.

where

$$d_1 = \frac{\ln(S/K) + (r + \sigma^2/2)(T-t)}{\sigma\sqrt{T-t}} \quad (15.16b)$$

$$d_2 = \frac{\ln(S/K) + (r - \sigma^2/2)(T-t)}{\sigma\sqrt{T-t}} = d_1 - \sigma\sqrt{T-t} \quad (15.16c)$$

and where  $N(x)$  denotes the standard cumulative normal probability distribution.

Let us examine some special cases. First suppose  $t = T$  (meaning the option is at expiration). Then

$$d_1 = d_2 = \begin{cases} +\infty & \text{if } S > K \\ -\infty & \text{if } S < K \end{cases}$$

because the  $d$ 's depend only on the sign of  $\ln(S/K)$ . Therefore, since  $N(\infty) = 1$  and  $N(-\infty) = 0$ , we find

$$C(S, T) = \begin{cases} S - K & \text{if } S > K \\ 0 & \text{if } S < K, \end{cases}$$

which agrees with the known value at  $T$ .

Next let us consider  $T = \infty$ . Then  $d_1 = \infty$  and  $e^{-r(T-t)} = 0$ . Thus  $C(S, \infty) = S$ , which agrees with the result derived earlier for a perpetual call.

**Example 15.2 (A 5-month option)** Let us calculate the value of the same option considered in Chapter 14, Example 14.3. That was a 5-month call option on a stock

with a current price of \$62 and volatility of 20% per year. The strike price is \$60 and the interest rate is 10%. Using  $S = 62$ ,  $K = 60$ ,  $\sigma = .20$ , and  $r = .10$ , we find

$$d_1 = \frac{\ln(62/60) + .12 \times 5/12}{.20\sqrt{5/12}} = .641287$$

$$d_2 = d_1 - .2\sqrt{5/12} = .512188.$$

The corresponding values for the cumulative normal distribution are found by the approximation in Exercise 1 to be

$$N(d_1) = .739332, \quad N(d_2) = .695740.$$

Hence the value for the call option is

$$C = 62 \times .739332 - 60 \times .95918 \times .695740 = \$5.798.$$

This is close to the value of \$5.85 found by the binomial lattice method.

Although a formula exists for a call option on a non-dividend-paying stock, analogous formulas do not generally exist for other options, including an American put option. The Black–Scholes equation, incorporating the corresponding boundary conditions, cannot be solved in analytic form.

## 15.4 Risk-Neutral Valuation\*

In the binomial lattice framework, pricing of options and other derivatives was expressed concisely as discounted risk-neutral valuation. This concept works in the Ito process framework as well.

For the geometric Brownian motion stock price process

$$dS(t) = \mu S dt + \sigma S dz \tag{15.17}$$

we know from Section 13.7 that

$$E[S(t)] = S(0)e^{\mu t}. \tag{15.18}$$

In a risk-neutral setting, the price of the stock at time zero is found from its price at time  $t$  by discounting the risk-neutral expected value at the risk-free rate. This means that there should hold

$$S(0) = e^{-rt}\hat{E}[S(t)].$$

It is clear that this formula would hold if  $\hat{E}[S(t)] = S(0)e^{rt}$ . From (15.17) and (15.18) this will be the case if we define the process

$$dS = rS dt + \sigma S d\hat{z}, \tag{15.19}$$

where  $\hat{z}$  is a standardized Wiener process, and we define  $\hat{E}$  as expectation with respect to the  $\hat{z}$  process. In other words, starting with a lognormal Ito process with rate  $\mu$ ,

we obtain the equivalent risk-neutral process by constructing a similar process but having rate  $r$ .

This change of equation is analogous to having two binomial lattices for a stock process: a lattice for the real process and a lattice for the risk-neutral process. In the first lattice the probabilities of moving up or down are  $p$  and  $1 - p$ , respectively. The risk-neutral lattice has the same values as the stock prices on the nodes, but the probabilities of up and down are changed to  $q$  and  $1 - q$ . For the Ito process we have two processes—like two lattices. Because the probability structures are different, we use  $z$  and  $\hat{z}$  to distinguish them.

Once the risk-neutral probability structure is defined, we can use risk-neutral valuation to value any security that is a derivative of  $S$ . In particular, for a call option the pricing formula is

$$C = e^{-rT} \hat{E}\{\max[S(T) - K, 0]\}. \quad (15.20)$$

This is analogous to equation (14.13) in Chapter 14.

We know that the risk-neutral distribution of  $S(T)$  satisfying (15.19) is lognormal with  $E\{\ln[S(T)/S(0)]\} = rT - \frac{1}{2}\sigma^2 T$  and  $\text{var}\{\ln[S(T)/S(0)]\} = \sigma^2 T$ . We can use this distribution to find the indicated expected value in analytic form. The result will be identical to the value given by the Black–Scholes equation for a call option price. Specifically, writing out the details of the lognormal distribution, we have

$$C = \frac{e^{-rT}}{\sqrt{2\pi\sigma^2 T}} \int_{\ln K}^{\infty} (e^x - K) e^x e^{-(x - \ln S(0) - rT + \sigma^2 T/2)^2/(2\sigma^2 T)} dx. \quad (15.21)$$

This is the Black–Scholes formula in integral form.

## 15.5 Delta

At any fixed time the value of a derivative security is a function of the underlying asset's price. The sensitivity of this function to changes in the price of the underlying asset is described by the quantity **delta** ( $\Delta$ ). If the derivative security's value is  $f(S, t)$ , then formally delta is

$$\Delta = \frac{\partial f(S, t)}{\partial S}.$$

Delta is frequently expressed in approximation form as

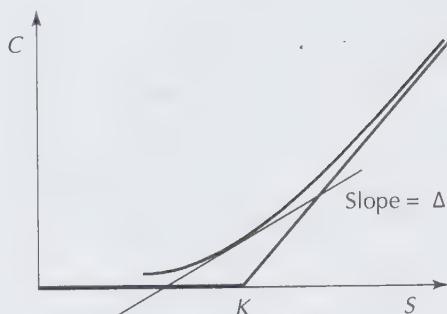
$$\Delta = \frac{\Delta f}{\Delta S}.$$

The delta of a call option is illustrated in Figure 15.2. It is the slope of the curve that relates the option price to the stock price.

Delta can be used to construct portfolios that hedge against risk. As an example, suppose that an option trader believes that a certain call option is overpriced. The trader would like to write (that is, sell) the option, taking a very large (negative) position in the call option. However, doing so would expose the trader to a great deal of price risk. If the underlying stock price should increase, the trader will lose money on the

**FIGURE 15.2 Delta of a call option.**

Delta measures the sensitivity of the option value to small changes in the price of the underlying security.



option even if his assessment of the option value relative to its current price is well founded. The trader may not wish to speculate on the stock itself, but only to profit from his belief that the option is overpriced. The trader can neutralize the effect of stock price fluctuations by offsetting the sale of options with a simultaneous purchase of the stock itself. The appropriate amount of stock to purchase is delta times the value of the options sold. Then if the stock price should rise by \$1, the profit on the trader's holding of stock will offset the loss on the options.

The delta of a call option can be calculated from the Black–Scholes formula (15.3) to be

$$\Delta = N(d_1). \quad (15.22)$$

This explicit formula can be used to implement delta hedging strategies that employ call options.

In general, given a portfolio of securities, all components of which are derivative to a common underlying asset, we can calculate the **portfolio delta** as the sum of the deltas of each component of the portfolio. Traders who do not wish to speculate on the underlying asset prices will form a portfolio that is **delta neutral**, which means that the overall delta is zero. In the case of the previous trader, the value of the portfolio was  $-C + \Delta \times S$ . Since the delta of  $S$  is 1, the overall delta of this hedged portfolio is  $-\Delta + \Delta = 0$ .

Delta itself varies both with  $S$  and with  $t$ . Hence a portfolio that is delta neutral initially will not remain so. It is necessary, therefore, **to rebalance** the portfolio by changing the proportions of its securities in order to maintain neutrality. This process constitutes a **dynamic hedging strategy**. In theory, rebalancing should occur continuously, although in practice it is undertaken only periodically or when delta has materially changed from zero.

The amount of rebalancing required is related to another constant termed **gamma** ( $\Gamma$ ). Gamma is defined as

$$\Gamma = \frac{\partial^2 f(S, t)}{\partial S^2}.$$

Gamma defines the curvature of the derivative price curve. In Figure 15.2 gamma is the second derivative of the option price curve at the point under consideration. For

a call option,

$$\Gamma = \frac{N'(d_1)}{S\sigma\sqrt{T}}.$$

Another useful number is **theta** ( $\Theta$ ). Theta is defined as

$$\Theta = \frac{\partial f(S,t)}{\partial t}.$$

Theta measures the time change in the value of a derivative security. Referring again to Figure 15.2, if time is increased, the option curve will shift to the right. Theta measures the magnitude of this shift. For a call option,

$$\Theta = -\frac{SN'(d_1)\sigma}{2\sqrt{T}} - rKe^{-rT}N(d_2).$$

These parameters are sufficient to estimate the change in value of a derivative security over small time periods, and hence they can be used to define appropriate hedging strategies. In particular, using  $\delta f$ ,  $\delta S$ , and  $\delta t$  to represent small changes in  $f$ ,  $S$ , and  $t$ , we have

$$\delta f \approx \Delta \cdot \delta S + \frac{1}{2}\Gamma \times (\delta S)^2 + \Theta \times \delta t$$

as a first-order approximation to  $\delta f$ .<sup>2</sup>

**Example 15.3 (Call price estimation)** Consider a call option with  $S = 43$ ,  $K = 40$ ,  $\sigma = .20$ ,  $r = 10\%$ , and a time to expiration of  $T - t = 6$  months = .5. The Black–Scholes formula yields  $C = \$5.56$ . We can also calculate that  $\Delta = .825$ ,  $\Gamma = .4235$ , and  $\Theta = -4.558$ . (See Exercise 7.)

Now suppose that in two weeks the stock price increases to \$44. We have  $\delta S = 1$  and  $\delta t = 1/26$ ; therefore the price of the call at that time is approximately

$$C \approx 5.56 + \Delta \times 1 + \frac{1}{2}\Gamma \times (1)^2 + \Theta \times (1/26) = \$6.233.$$

The actual value of the call at the later date according to the Black–Scholes formula is  $C = \$6.23$ .

## 15.6 Replication, Synthetic Options, and Portfolio Insurance\*

The derivation of the Black–Scholes equation shows that a derivative security can be duplicated by constructing a portfolio consisting of an appropriate combination of the underlying security and the risk-free asset. We say that this portfolio **replicates** the derivative security. The proportions of stock and the risk-free asset in the portfolio must be adjusted continuously with time, but no additional money need be added or taken away; the portfolio is self-financing. This replication can be carried out in practice in order to construct a **synthetic** derivative security using the underlying

<sup>2</sup> Recall that  $\delta S$  is proportional to  $\sqrt{\delta t}$ , so we must include the  $(\delta S)^2$  term.

and the risk-free assets. Of course, the required construction is dynamic, since the particular combination must change every period (or continuously in the context of the Black–Scholes framework).

The process for a call option is this: At the initial time, calculate the theoretical price  $C$ . Devote an amount  $C$  to the replicating portfolio. This portfolio should have  $\Delta S$  invested in the stock and the remainder invested in the risk-free asset (although this will usually require *borrowing*, not lending). Then both the delta and the value of the portfolio will match those of the option. Indeed, the short-term behavior of the two will match.

A short time later, delta will be different, and the portfolio must be rebalanced. However, the value of the portfolio will be approximately equal to the corresponding new value of the option, so it will be possible to continue to hold the equivalent of one option. This rebalancing is repeated frequently. As the expiration date of the (synthetic) option approaches, the portfolio will consist mainly of stock if the price of the stock is above  $K$ ; otherwise the portfolio's value will tend to zero.

**Example 15.4 (A replication experiment)** Let us construct, experimentally, a synthetic call option on Exxon stock with a strike price of \$35 and a life of 20 weeks. We will replicate this option by buying Exxon stock and selling (that is, borrowing) the risk-free asset. We select a 20-week period. The weekly closing prices are shown in the second column of Table 15.1. The measured sigma corresponding to this period is  $\sigma = 18\%$  on an annual basis, so we shall use that value to calculate the theoretical values of call prices and delta. We assume an interest rate of 10%.

Let us walk across the first row of the table. There are 20 weeks remaining in the life of the option. The initial stock price is \$35.50. The third column shows that the initial value of the call (as calculated by the Black–Scholes formula) is \$2.62. Likewise the initial value of delta is .701. To construct the replicating portfolio we devote a value of \$2.62 to it, matching the initial value of the call. This is shown in the column marked “Portfolio value.” However, this portfolio consists of two parts, indicated in the next two columns. The amount devoted to Exxon stock is \$24.89, which is delta times the current stock value. The remainder  $\$2.62 - \$24.89 = -\$22.27$  is devoted to the risk-free asset. In other words we borrow \$22.27, add \$2.62, and use the total of \$24.89 to buy Exxon stock.

Now walk across the second row, which is calculated in a slightly different way. The first four entries show that there are 19 weeks remaining, the new stock price is \$34.63, the corresponding Black–Scholes option price is \$1.96, and delta is now .615. The next entry, “Portfolio value,” is obtained by updating from the row above it. The earlier stock purchase of \$24.89 is now worth  $(34.63/35.50) \times \$24.89 = \$24.28$ . The debt of \$22.27 is now a debt of  $(1 + 0.10/52)\$22.27 = \$22.31$ . The new value of the portfolio we constructed last week is therefore now  $\$24.28 - \$22.31 = \$1.96$  (adjusting for the round-off error in the table). This new value does not exactly agree with the current call value (although in this case it happens to agree within the two decimal places shown). We do not add or subtract from the value. However, we now rebalance the portfolio by allocating to the stock \$21.28 (which is delta times the stock price) and borrowing \$19.32 so that the net portfolio value remains at \$1.96.

**TABLE 15.1**  
**AN EXPERIMENT IN OPTION REPLICATION**

Weeks remaining	XON price	Call price	Delta	Portfolio value	Stock portfolio	Bond portfolio
20	35.50	2.62	.701	2.62	24.89	-22.27
19	34.63	1.96	.615	1.96	21.28	-19.32
18	33.75	1.40	.515	1.39	17.37	-15.98
17	34.75	1.89	.618	1.87	21.47	-19.59
16	33.75	1.25	.498	1.22	16.79	-15.58
15	33.00	0.85	.397	.81	13.09	-12.28
14	33.88	1.17	.494	1.14	16.74	-15.60
13	34.50	1.42	.565	1.41	19.48	-18.07
12	33.75	0.96	.456	.96	15.39	-14.43
11	34.75	1.40	.583	1.38	20.27	-18.89
10	34.38	1.10	.522	1.13	17.94	-16.81
9	35.13	1.44	.624	1.49	21.92	-20.43
8	36.00	1.94	.743	2.00	26.74	-24.75
7	37.00	2.65	.860	2.69	31.80	-29.11
6	36.88	2.44	.858	2.53	31.65	-29.12
5	38.75	4.10	.979	4.08	37.92	-33.84
4	37.88	3.17	.961	3.16	36.39	-33.23
3	38.00	3.21	.980	3.22	37.25	-34.03
2	38.63	3.76	.998	3.76	38.56	-34.79
1	38.50	3.57	1.000	3.57	38.50	-34.93
0	37.50	2.50		2.50		

A call on XON with strike price 35 and 20 weeks to expiration is replicated by buying XON stock and selling the risk-free asset at 10%. The portfolio is adjusted each week according to the value of delta at that time. When the volatility is set at 18% (the actual value during that period), the portfolio value closely matches the Black-Scholes value of the call.

Succeeding rows are calculated in the same fashion. At each step, the updated portfolio value may not exactly match the current value of the call, but it tends to be very close, as is seen by scanning down the table and comparing the call and portfolio values. The maximum difference is 11 cents. At the end of the 20 weeks it happens in this case that the portfolio value is exactly equal (to within a fraction of a cent) to the value of the call.

The results depend on the assumed value of volatility. The choice of  $\sigma = 18\%$  represents the actual volatility over the 20-week period, and this choice leads to good results. Study of a longer period of Exxon stock data before the date of this option indicates that volatility is more typically 20%. If this value were used to construct Table 15.1, the resulting final portfolio value would be \$2.66 rather than \$2.50. If  $\sigma = 15\%$  were used, the final portfolio value would be \$2.27.

The degree of match would also be affected by transactions costs. The experiment with an Exxon call assumed that transactions costs were zero and that stock could be purchased in any fractional amount. In practice these assumptions are not satisfied exactly. But for large volumes, as might be typical of institutional dealings,

the departure from these assumptions is small enough so that replication is in fact practical.

**Example 15.5 (Portfolio insurance)** Many institutions with large portfolios of equities (stocks) are interested in insuring against the risk of a major market downturn. They could protect the value of their portfolio if they could buy a put, giving them the right to sell their portfolio at a specified exercise price  $K$ .

Puts are available for the major indices, such as the S&P 500, and hence one way to obtain protection is to buy index puts. However, a particular portfolio may not match an index closely, and hence the protection would be imperfect.

Another approach is to construct a synthetic put using the actual stocks in the portfolio and the risk-free asset. Since puts have negative deltas, construction of a put requires a short position in stock and a long position in the risk-free asset. Hence some of the portfolio would be sold and later bought back if the market moves upward. This strategy has the disadvantage of disrupting the portfolio and incurring trading costs.

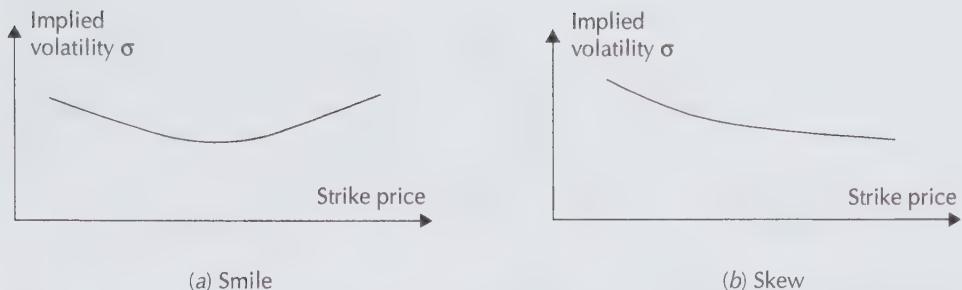
A third approach is to construct a synthetic put using futures on the stocks held in the portfolio instead of using the stocks themselves. To implement this strategy, one would calculate the total value of the puts required and go long delta times this amount of futures. (Since  $\Delta < 0$ , we would actually short futures.) The difference between the value of stock shorted and the value of a put is placed in the risk-free asset. The positions must be adjusted periodically as delta changes, just as in the previous example. This method, termed **portfolio insurance**, was quite popular with investment institutions (such as pension funds) for a short time until the U.S. stock market fell substantially in October 1987, and it was not possible to sell futures in the quantities called for by the hedging rule, resulting in loss of protection and actual losses in portfolio value.

## 15.7 Volatility Smiles

Traders often state their opinion about the value of an option, not directly in terms of a price, but in terms of the volatility they would use to calculate the price with the Black–Scholes equation. After all, to compute the Black–Scholes price, one must provide the risk-free rate, the strike price, the maturity date, and, finally, the volatility. All of these parameters except volatility are basically objective; all traders should agree on their values. The differences in the assessed value of the option boil down to differences in volatility estimates. Conversely, once a price is stated, it is possible to reverse the Black–Scholes formula and determine the volatility that would produce that price. This is said to be the **implied volatility** of the option.

To this point, the notion of implied volatility is merely a statement about the relation between independent and dependent variables and how they can be reversed. But the notion turns out to be a bit more interesting; further pursuit of the subject provides a lesson that may deepen our understanding of modern pricing theory.

We start with the observation that various options on the same asset often do not have the same implied volatility, contrary to what would hold for the Black–Scholes



**FIGURE 15.3 Implied volatility curves for two cases.**

equation based on an underlying geometric Brownian motion process. By definition, all option prices on the same underlying asset would have the same implied volatility. That is in theory. In practice, the implied volatilities vary somewhat, but in a systematic way.

Consider a group of foreign exchange call options on the same asset, with the same maturity but with different strike prices. A typical curve of implied volatility versus strike price is shown in Figure 15.3(a). It has a characteristic **smile** shape. A different shape, shown in Figure 15.3(b), is typical for equity options. Its shape is referred to as a **skew**. Other common shapes have suggestive names such as “smirk” and “sneer.”

## Equality of Implied Volatilities

It can be argued that the implied volatility of a put is always equal to the implied volatility of the corresponding call (the call with the same strike price and maturity date as the put). This follows from put–call parity. The standard argument for parity is based on an arbitrage argument and hence applies to market prices. Theoretically it must also follow for Black–Scholes prices. Therefore, from Section 14.3 we may write the put–call parity equation, first for the market and second for Black–Scholes prices,

$$C_M - P_M = S - dK$$

$$C_{BS} - P_{BS} = S - dK,$$

where  $C_M$  and  $P_M$  are the market prices of the call and put and  $C_{BS}$  and  $P_{BS}$  are the corresponding Black–Scholes prices. This leads directly to

$$C_M - C_{BS} = P_M - P_{BS}. \quad (15.23)$$

Now suppose that the assumed volatility is adjusted to make the Black–Scholes call price equal to the market call price, forcing the left side of equation (15.23) to zero.

The right side must also go to zero, which means that the implied volatility of the put must equal that of the call. In other words, the implied volatility curve (the smile) of the put must be identical to that of the call. From this observation it makes sense to define an implied volatility function  $\sigma(K, T)$ , a function of strike price  $K$  and maturity date  $T$ , that applies to all European options, call or put, on a given asset. Of course in general, we may have direct measurement of this function at only a few combinations of  $K$  and  $T$ .

## Risk-Neutral Probability Density\*

The volatility smile is an interesting phenomenon with significant practical value if someone wishes to trade options or (as we shall see) other, more complex derivatives. However, it also provides an opportunity to investigate further the role and implications of risk-neutral pricing. Basically, we trace the source of the smile phenomenon back to the underlying probability structure.

If the price of an asset follows geometric Brownian motion, we know that the probability density of the price at time  $T$  is lognormal, with log mean  $vT$  and log variance  $\sigma^2 T$ . Furthermore, we know that the corresponding risk-neutral density is also lognormal but with log mean  $rT$  and log variance  $\sigma^2 T$ . This knowledge leads, in fact, to the explicit computation of the price of an option as the discounted risk-neutral expected value of the option's payoff—as described in equation (15.21). We know, however, that the assumption of a lognormal density for actual prices is only an approximation, with the true density likely having heavy tails and/or skew characteristics. Clearly there is no reason to expect that the corresponding risk-neutral density will be lognormal. Rather, we expect that it will inherit, in some perhaps-distorted manner, the heavy tails and/or skewness of the real probability density.

Let us define the risk-neutral density as  $g(S, T)$ . Then the current value of a payoff  $F(S, T)$  at time  $T$  is found from the discounted risk-neutral expectation

$$V = e^{-rT} \int_{-\infty}^{\infty} F(S, T) g(S, T) dS.$$

If we let  $C(K, T)$  denote the value of a European call option with strike  $K$  and maturity  $T$ , we have

$$C(K, T) = e^{-rT} \int_K^{\infty} (S - K) g(S, T) dS.$$

Differentiating with respect to  $K$  we find

$$\frac{\partial C(K, T)}{\partial K} = -e^{-rT} \int_K^{\infty} g(S, T) dS.$$

And differentiating again,

$$\frac{\partial^2 C(K, T)}{\partial K^2} = e^{-rT} g(K, T).$$

Reversing this last equation we find

$$g(K, T) = e^{rT} \frac{\partial^2 C(K, T)}{\partial K^2}, \quad (15.24)$$

which shows that the density can be recovered from the call price function.

Indeed for fixed  $T$ , the price function  $C(K, T)$ , the risk-neutral density function  $g(S, T)$ , and the implied volatility function  $\sigma(K, T)$  all contain the same information. Each can be found from any one of the others. For example,  $C(K, T)$  and  $\sigma(K, T)$  are equivalent by use of the Black–Scholes formula. Also,  $g(K, T)$  can be found from equation (15.24); and  $C(K, T)$  can be found from  $g(K, T)$  by evaluating the discounted risk-neutral expected value of calls. Of course some of these transformations are complex, requiring the inversion (that is, explicit solution) of differential equations. Special numerical techniques have been developed to find good approximations.

The maturity date  $T$  can also be varied, and then  $\sigma(K, T)$  as a function of  $T$  is the **term structure of volatility**. As both  $K$  and  $T$  vary, there is a **volatility surface**. These functions provide a means for estimating the risk-neutral density function  $g(S, T)$ .

Let us apply what we have studied to explain the shape of the implied volatility curve. Consider an asset with a probability density that has heavy tails. We expect that the risk-neutral density will also have this characteristic. Now consider a European call option on this asset that is well out of the money. The price of the call will be largely dependent on the risk-neutral probability in the upper tail because a large increase in the price of the underlying asset is necessary to get the call into the money. If the risk-neutral upper tail is heavier than the tail of a lognormal density, we expect that the price of the option will be greater than the Black–Scholes price, and, consequently, the implied volatility will be relatively higher than the lognormal volatility. A similar argument shows that an out-of-the-money European put will have a relatively high implied volatility if the actual density has a heavy lower tail. The two can be combined to produce the smile shape, since, according to equation (15.23), the implied volatilities for puts and calls are the same.

For an asset with a downwardly skewed density, we expect that the implied volatility will decrease with the strike price, and hence the curve of implied volatility will take the skew form. These predictions are well supported by observation.

We see now that the volatility smile conveys information about the risk-neutral density. By observing the prices of a sufficiently large number of options with different strike prices and different maturity dates, we can, by inverting the relations discussed in this section, form a reasonable approximation to the risk-neutral density  $g(S, T)$ , and this can provide a basis for evaluating other, perhaps more complex, derivatives of the underlying asset.

## 15.8 Computational Methods

The theory presented in this chapter can be transformed into computational methods in several ways. Some of these methods are briefly outlined in this section.

## Monte Carlo Simulation

**Monte Carlo simulation** is one of the most powerful and most easily implemented methods for the calculation of option values. However, the procedure is essentially only useful for European-style options, where no decisions are made until expiration. Suppose that there is a derivative security that has payoff at the terminal time  $T$  of  $f(S(T))$  and suppose the stock price  $S(t)$  is governed by geometric Brownian motion according to

$$dS = \mu S dt + \sigma S dz,$$

where  $z$  is a standardized Wiener process. The basis for the Monte Carlo method is the risk-neutral pricing formula, which states that the initial price of the derivative security should be

$$P = e^{-rT} \hat{E}[f(S(T))].$$

To evaluate the right-hand side by Monte Carlo simulation, the stochastic stock dynamic equation in a risk-free world

$$dS = rS dt + \sigma S d\hat{z}$$

is simulated over the time interval  $[0, T]$  by dividing the entire time period into several periods of length  $\Delta t$ . The simulation equation is

$$S(t_k + \Delta t) = S(t_k) + rS(t_k)\Delta t + \sigma S(t_k)\epsilon(t_k)\sqrt{\Delta t},$$

where  $\epsilon(t_k)$  is chosen by a random number generator that produces numbers according to a normal distribution having zero mean and variance 1. (Or the multiplicative version of Section 13.7 can be used.) After each simulation, the value  $f(S(T))$  is calculated. An estimate  $\hat{P}$  of the true theoretical price of the derivative security is found from the formula

$$\hat{P} = e^{-rT} \text{average}[f(S(T))],$$

where the average is taken over all simulation trials.

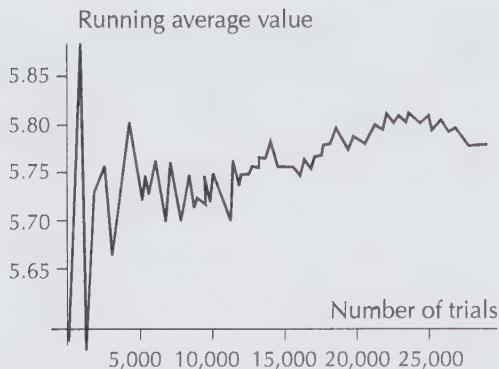
A disadvantage of this method is that suitable accuracy may require a very large number of simulation trials. In general, the expected error decreases with the number of trials  $n$  by the factor  $1/\sqrt{n}$ ; so one more digit of accuracy requires 100 times as many trials. Often tens of thousands of trials are required to obtain two-place accuracy.

**Example 15.6 (The 5-month call)** Simulation is unnecessary for a call option since better methods are available, but this example, which was solved earlier in Example 15.2, provides a simple illustration of the method. For this call  $S(0) = \$62$ ,  $K = \$60$ ,  $\sigma = 20\%$ , and  $r = 12\%$ . The time to maturity is 5 months.

To carry out the simulation the 5-month period was divided into 80 equal small time intervals. The stock dynamics were modeled as

$$S(t + \Delta t) = S(t) + rS(t)\Delta t + \sigma S(t)\epsilon(t)\sqrt{\Delta t},$$

where  $\epsilon(t)$  is chosen randomly from a normal distribution with mean zero and unit variance.



**FIGURE 15.4 Monte Carlo evaluation of a call.** The value of a call is estimated as the discounted average of final payoff when simulations are governed by the risk-neutral process. The method is easy to implement but requires a large number of trials for reasonable accuracy.

After each simulation trial, the terminal value of the call,  $\max(S - K, 0)$ , was determined based on the final stock price, and this value was discounted back to the initial time. A running average of these discounted values was recorded as successive runs were made. Figure 15.4 shows a graph of the discounted average value obtained as a function of the total number of trials. A reasonably accurate and stable result requires about 25,000 simulation trials. From the figure we can conclude that the price of the call is in the neighborhood of \$5.80 plus or minus around 10 cents. The Black–Scholes value is in fact \$5.80.

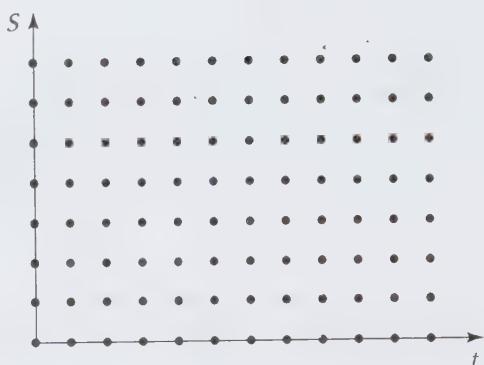
The simulation can be improved by various **variance reduction** procedures, the two most common of these being the **control variate method** and the **antithetic variable method**. (See Exercise 13.)

Although it is costly in terms of computer time to use the Monte Carlo method, the method is in fact often used in practice to evaluate European-style derivatives that do not have analytic solutions. The method has the advantages of flexibility and ease of programming, and it is reasonably foolproof.

## Finite-Difference Methods

Numerical solution of the Black–Scholes partial differential equation is a second approach to the calculation of option prices. In this method a large rectangular grid is established, a small version of which is shown in Figure 15.5. In this grid the horizontal axis represents time  $t$  and the vertical axis represents  $S$ . The time difference between horizontally adjacent points is  $\Delta t$ , and the price difference between vertically adjacent points is  $\Delta S$ . The function  $f(S, t)$  is defined at all the corresponding grid points. If the  $S$  values on the grid are indexed by  $i$  and the  $t$  values are indexed by  $j$ , then the function at the grid point  $(i, j)$  is denoted by  $f_{ij}$ .

**FIGURE 15.5 Grid for finite-difference method.** The finite-difference method approximates the Black–Scholes equation by algebraic relations among values at grid points. The method can handle American as well as European options.



The method is implemented by using the finite-difference approximations to partial derivatives as follows:

$$\frac{\partial f}{\partial S} \approx \frac{f_{i+1,j} - f_{i,j}}{\Delta S}$$

$$\frac{\partial^2 f}{\partial S^2} \approx \frac{f_{i+1,j} - f_{i,j} - f_{i-1,j} + f_{i-1,j}}{(\Delta S)^2} = \frac{f_{i+1,j} - 2f_{i,j} + f_{i-1,j}}{(\Delta S)^2}$$

$$\frac{\partial f}{\partial t} \approx \frac{f_{i,j+1} - f_{i,j}}{\Delta t}.$$

The terminal conditions imply that  $f_{i,j}$  is known at the right boundary of the grid. Additional boundary conditions may be specified, depending on the particular derivative security. In the case of a put option, for example, it is known that the value of the put is at least equal to  $K - S$  everywhere, and since the value of the put approaches zero as  $S \rightarrow \infty$ , we may specify that the value is zero along the top edge of the grid.

When these approximations are used in the Black–Scholes equation, the result is a large set of algebraic equations and inequalities. These can be solved systematically by working backward from the right edge of the grid toward the left. In fact, the equations are closely related to the equations of backward solution in a lattice.

The finite-difference method has the advantage that it can handle derivative securities such as American puts that impose boundary conditions other than terminal-time conditions. An inherent disadvantage, however, is that the equations are only approximations to the actual partial differential equation, and therefore, aside from the obvious approximation error, their solutions are subject to instabilities and inconsistencies, which are not characteristic of the partial differential equation itself (usually resulting from implied probabilities becoming negative). As a general rule of numerical problem solving, if a problem is to be solved with a finite-step approximation, it is usually better to reformulate the problem itself in finite-step form and then solve that problem directly, rather than to formulate the problem in continuous time and then approximate the solution by a finite-step method. In the case of derivative securities this means that rather than approximating the Black–Scholes equation, it is probably

better to use a discrete formulation, such as the discrete-time risk-neutral pricing formula or the binomial lattice formulation. These discrete formulations will introduce approximation error, but will not instill numerical instabilities. Despite these caveats, finite-difference methods, when carefully designed, do have a useful role in the numerical evaluation of derivative securities.

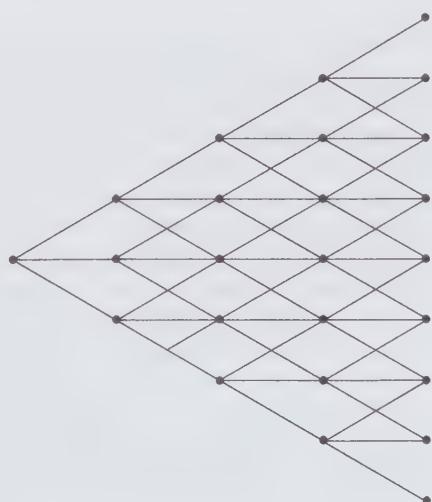
## Binomial and Trinomial Lattices

A popular method for finding the value of a derivative security is the binomial lattice method of Section 14.6. The method is straightforward and leads to reasonably accurate results, even if the time divisions are crude (say, 10 or so time periods over the remaining time interval). However, it is also possible to use other tree and lattice structures. For example, a good choice is to use a trinomial lattice, as shown in Figure 15.6. For a given number of time periods, the trinomial lattice has more nodes than a binomial lattice and hence can produce a better approximation to the continuous solution.

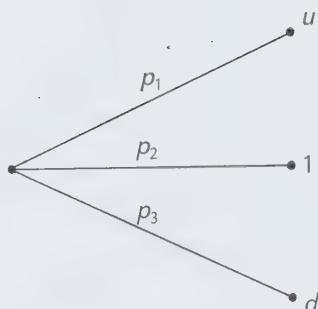
At first it might seem that a trinomial lattice cannot replace a binomial lattice because it is impossible to replicate three possible outcomes using only two securities: the stock and the risk-free asset. This is correct; replication is not possible. Hence the trinomial lattice cannot be used as a basis for options theory. However, once the theory is deduced by other methods (such as the Black–Scholes method), we can seek alternative ways to implement it. A trinomial lattice is a convenient structure for implementing the risk-neutral pricing formula.

To set up a suitable trinomial lattice refer to Figure 15.7, which shows one piece of the lattice. There are three paths leaving a node, with probabilities  $p_1$ ,  $p_2$ , and  $p_3$ . The three resulting nodes represent multiplication of the stock value by  $u$ , 1, and  $d$ , respectively, where we set  $d = 1/u$ , so that an up followed by a down is equal to 1.

**FIGURE 15.6 Trinomial lattice.** A trinomial lattice can give a more accurate representation than a binomial lattice for the same number of steps.



**FIGURE 15.7 One piece of a trinomial lattice.** In this lattice we must have  $d = 1/u$  so that the nodes recombine after two steps.



To assign the parameters of the trinomial lattice we can arbitrarily select a value for  $u$ . Then if the mean value for one step is to be  $1 + \mu\Delta t$  and the variance is to be  $\sigma^2\Delta t$ , we select the probabilities to satisfy

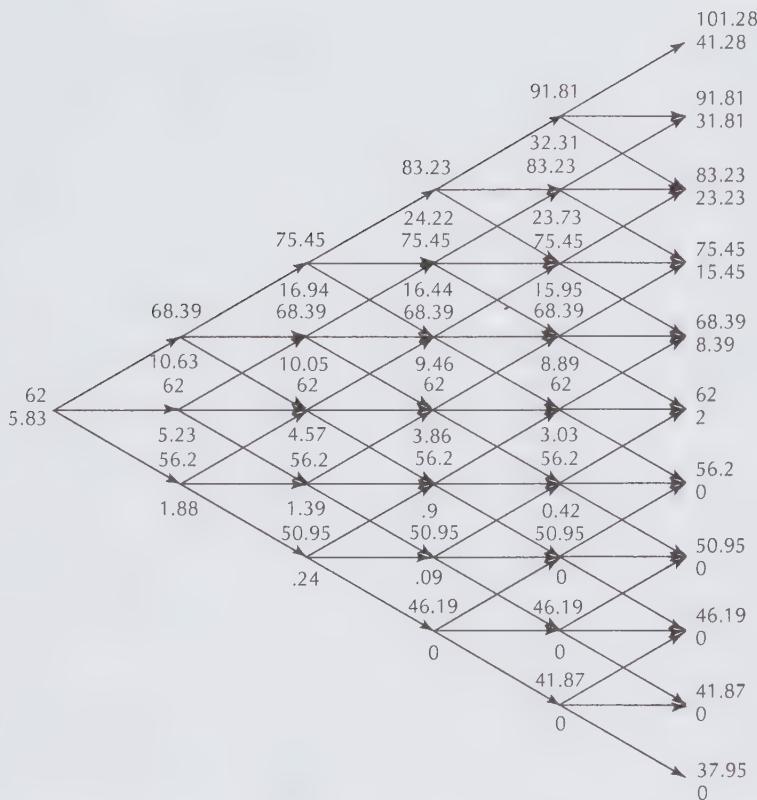
$$\begin{aligned} p_1 + p_2 + p_3 &= 1 \\ up_1 + p_2 + dp_3 &= 1 + \mu\Delta t \\ u^2p_1 + p_2 + d^2p_3 &= \sigma^2\Delta t + (1 + \mu\Delta t)^2. \end{aligned} \tag{15.25}$$

(The last line represents  $E(x^2) = \text{var}(x) + E(x)^2$ , where  $x$  is the random factor by which the stock price is multiplied in one period.) This is just a system of three linear equations to be solved for the three probabilities. Once these probabilities are found, we have a good approximation to the underlying stock dynamics. (Note that we are implicitly using the dynamics of (13.19).)

To use this lattice for pricing, we must instead use the risk-neutral probabilities  $q_1$ ,  $q_2$ , and  $q_3$ . These are found by solving the same set of equations (15.25), but with the mean value changed from  $\mu\Delta t$  to  $r\Delta t$ . Once the risk-neutral probabilities are found, the lattice can be solved backward, just as in the binomial procedure.

**Example 15.7 (The 5-month call)** Let us find the price of the 5-month call option of Example 14.3 using a trinomial lattice, just to compare the results. We have  $S(0) = \$62$ ,  $K = \$60$ ,  $r = 10\%$ , and  $\sigma = 20\%$ . The time to expiration is 5 months = .416667. To set up the lattice we must select a value of  $u$  and solve the equations (15.25) for the probabilities (when  $\mu$  is set to  $r$ ) in the equations. The choice of  $u$  requires a bit of experimentation, since for some values the resulting risk-neutral probabilities may not be positive. For example, using  $u = 1.06$  leads to  $q_1 = .57$ ,  $q_2 = -.03$ , and  $q_3 = .46$ . Instead we use  $u = 1.1031277$  and  $q_1 = .20947$ ,  $q_2 = .64896$ , and  $q_3 = .14156$ . This leads to the lattice shown in Figure 15.8. Note that the value of the option obtained is \$5.83, which is slightly closer to the Black–Scholes result of \$5.80 than is the price of \$5.85 determined by a binomial lattice.<sup>3</sup>

<sup>3</sup> In this example we assumed monthly compounding, while the Black–Scholes formula implicitly assumes continuous compounding. We can also use the equivalent continuous compounding rate in the example, and the result differs by only one-tenth of a cent from \$5.83.



**FIGURE 15.8 5-month call using a trinomial lattice.** Stock prices are listed above nodes; and option prices are listed below. The discounted risk-neutral valuation is easily generalized to the trinomial lattice.

The lattice of Figure 15.8 has the stock value listed above each node and the option value listed below each node. The final option values are just  $\max(0, S - K)$ . The option values at other nodes are found by discounted risk-neutral pricing. For example, the value at the top node after 4 months is  $(1 + .10/12)^{-1}(q_1 \times 41.28 + q_2 \times 31.81 + q_3 \times 23.23) = 32.31$ . If in this calculation the stock values 101.28, 91.81, and 83.23 were used instead of the option values, the result would be the stock value of 91.81, but of course it is not necessary to use this backward procedure for the stock prices.

## 15.9 Exotic Options

Numerous variations on the basic design of options have been proposed. Each variation offers effective control of the risk perceived by a certain group of investors or eases execution and bookkeeping. We list a few of these variations here:

- 1. Bermudan option** In this option, the allowable exercise dates are restricted, in some cases to specific dates and in other cases to specific periods within the lifetime of the option. Warrants on stock often have this characteristic.
- 2. Forward start options** These are options that are paid for at one date, but do not begin until a later date.
- 3. Compound options** A compound option is an option on an option.
- 4. “As you like it” or “chooser” options** The holder of an “as you like it” option can, after a specified time, declare the option to be either a put or a call.
- 5. CAPs** These options restrict the amount of profit that can be made by the option holder by automatically exercising once the profit reaches a specified level. A \$20 CAP on a call option, means that once the stock price rises to \$20 over the strike price, the option is exercised.
- 6. LEAPS** This term stands for “Long-term Equity Anticipating Securities.” They are long-term, exchange-traded options with exercise dates as far as 3 years into the future.
- 7. Digital options** In a digital option the payoff is \$1 if the option is in the money and zero otherwise. A European digital call option, for example, has payoff 1 if  $S(T) > K$ , and 0 if  $S(T) < K$ , where  $K$  is the strike price.
- 8. Exchange options** Such an option gives one the right to exchange one specified security for another.
- 9. Yield-based options** A yield-based option on a bond defines the exercise value in terms of yield rather than price. Hence the holder of a yield-based call option benefits if bond prices **decrease** since yields move in the opposite direction to prices.
- 10. Cross-ratio options** These are foreign-currency options denominated in another foreign currency; for example, a call on German marks with an exercise price in Japanese yen.
- 11. Knockout options** These options terminate (with zero value) once the price of the underlying asset reaches a specified point. For calls these are “down and out” options, which terminate once the price of the underlying asset falls below a specified level. For puts the analogous option is a “up and out” option.
- 12. Discontinuous options** These options have payoffs that are discontinuous functions of the price of the underlying asset. For example, a call option may pay either zero or \$20, depending on whether the final price of the underlying asset is below or above a specified strike price.
- 13. Lookback options** In a lookback option the effective strike price is not specified, but is determined by the minimum (in the case of a call) or maximum (in the case of a put) of the price of the underlying asset during the period of the option. For example, a European-style lookback call option has a payoff equal to  $\max(S_T - S_{\min}, 0) = S_T - S_{\min}$ , where  $S_{\min}$  is the minimum value of the price  $S$  over the period from initiation to the termination time  $T$ . Such options are very attractive to investors, since in fact they always have positive value (unless  $S_T = S_{\min}$ ). Of course their prices reflect the apparent attractiveness.

**14. Asian options** The payoff of Asian options depends on the average price  $S_{\text{avg}}$  of the underlying asset during the period of the option. There are basically two ways that the average can be used. In one,  $S_{\text{avg}}$  serves as the strike price, so that the payoff of a corresponding call, for example, is  $\max(S_T - S_{\text{avg}}, 0)$ . In the second type,  $S_{\text{avg}}$  is substituted for the final price. Thus the payoff of the corresponding call is  $\max(S_{\text{avg}} - K, 0)$ , where  $K$  is a specified strike price.

## Pricing\*

Prices of some of these variations can be worked out computationally by using the theory and methods presented in this chapter. In other cases, formulas analogous to the Black–Scholes formula have been derived. There are cases, however, that present a serious technical challenge to the investment analysis community.

**Example 15.8 (A down and outer)** Consider a down and out call option on a non-dividend-paying stock. This option has a strike price of  $K$  and a “knockout” price of  $N < K$ . If the stock price  $S$  falls below  $N$ , the option is terminated with zero value. A closed-form expression for the original value of such an option can be found using the Black–Scholes framework; however, the details are not neat. We shall consider a simplified case, where the option is perpetual (that is,  $T = \infty$ ) but still has the down and out provision.

Since there is no explicit time dependence in the price of a perpetual option, the Black–Scholes equation reduces to

$$\frac{1}{2}\sigma^2 S^2 C''(S) + rSC'(S) - rC(S) = 0. \quad (15.26)$$

The boundary condition is

$$C(N) = 0.$$

We also know that  $C(S) \approx S$  as  $S \rightarrow \infty$ .

To solve (15.26) let us try a solution of the form  $C(S) = S^\alpha$ . This gives the algebraic equation

$$\frac{1}{2}\sigma^2 \alpha(\alpha - 1) + r\alpha - r = 0$$

which has solutions  $\alpha = 1$  and  $\alpha = -\gamma$ , where  $\gamma = 2r/\sigma^2$ . We may write the general solution of equation (15.26) as a linear combination of these two; that is,

$$C(S) = a_1 S + a_2 S^{-\gamma}.$$

Using the boundary condition we find  $a_2 = -a_1 N^{-\gamma}$ . Hence  $C(S) = a_1 [S - N(S/N)^{-\gamma}]$ . Using the asymptotic property, we have  $a_1 = 1$ . Therefore the final result is

$$C(S) = S - N(S/N)^{-\gamma}.$$

Since the value of a perpetual call is  $S$ , the second term in this expression can be regarded as a discount for the down and out feature.

The lookback and Asian options are particularly interesting because their payoffs are **path dependent**; that is, their payoffs do not merely depend on the final value

of the price of the underlying asset, but also on the way that that price was reached. So the conventional binomial lattice method of evaluation is not applicable. However, there are ways to modify the lattice approach to handle such cases; but as one might expect, the amount of computation required tends to be substantially greater than for a conventional option.

For European-style options that are path dependent, the Monte Carlo method offers a simple and effective procedure. The principle of risk-neutral pricing still applies, so it is only necessary to simulate the process repeatedly, using the risk-neutral probabilities for the underlying asset price fluctuations, and to average the payoffs obtained during the simulations. The control variate method for reducing the number of required simulations is especially useful for these options, and a corresponding non-path-dependent option, for which a solution is readily found, can be used as the control variate. (See Exercise 13.)

## 15.10 Comparison of Methods

Each computational procedure has its unique set of features, and these often guide the choice of which method to employ for a given derivative-pricing problem. Although it is difficult to make definitive general statements, Table 15.2 summarizes two of the main distinctions among methods; in particular, the ability to evaluate derivatives with decisions (such as when to exercise an option) and the ability to evaluate derivatives with values that are path dependent. These distinctions are not hard and fast, for sometimes a method can be augmented to handle cases that the standard version of the method cannot.

Lattice methods are used frequently in the text, because they are easy to program for small to modest-size problems. A major advantage is that they can optimize decisions, such as when to exercise an option or purchase new equipment. A disadvantage is that they cannot (easily) treat path-dependent derivatives.

Monte Carlo simulation is a basic tool of derivative evaluation. It is relatively easy to program even if the situation is complex (such as in a CMO). Its major drawbacks are that it cannot optimize ongoing decisions and it can be slow.

Trees get large exponentially with the number of time steps, but there are software packages that handle them fairly well. Trees have the great advantage that they can handle both decisions and path dependencies.

**TABLE 15.2  
METHOD FEATURES**

	Lattice	Simulation	Tree	Finite difference
Decisions	Y	N	Y	Y
Path dependent	N	Y	Y	N

The best computational method for pricing is often influenced by the nature of the problem.

Finite-difference methods provide natural ways to treat derivatives, since the general derivative theory is commonly cast in the form of partial differential equations. These methods discretize the variables and essentially convert the problem to a large set of equations, and these are solved by iterative methods. They can be solved backward, like a lattice, and then are able to treat decisions. Care must be taken to ensure convergence.

There are clever ways to minimize the disadvantages inherent in some methods. For example, consider an American-style put option, which has the possibility of early exercise. Using a Monte Carlo approach, one might hypothesize that there is a curve of the form  $S_E(t)$  that defines the exercise boundary; that is, if at a given time  $t$  we find  $S \leq S_E(t)$ , then the option is exercised at that point. A Monte Carlo or finite-difference method can evaluate the option with this prescribed exercise policy. Then various versions of the curve  $S_E(t)$  can be tried—perhaps by parameterizing it and using different parameter values—to find the best curve.

## 15.11 Storage Costs and Dividends\*

Commodity storage costs and security dividends can complicate an evaluation procedure, but there is an important special case, of proportional costs or dividends, that can be handled easily. This case is useful in applications, and the study of the technique involved should further enhance your understanding of risk-neutral pricing.

### Binomial Form

Suppose the commodity price  $S$  is governed by a binomial process having an up factor  $u$  and a down factor  $d$ . There is a storage cost of  $c$  per period, payable at the end of each period. The total risk-free return per period is  $R$ .

If you invest in the commodity at the beginning of a period, you must pay the current price  $S$ . At the end of the period, you receive the new commodity minus the storage cost; hence you receive either  $(u - c)S$  or  $(d - c)S$ . The new factors  $u - c$  and  $d - c$  are the legitimate factors that define the result of holding the commodity, and therefore these are the factors that can be used in a replication argument. It follows that the risk-neutral probabilities for up and down are

$$q = \frac{R - d + c}{u - d}, \quad 1 - q = \frac{u - c - R}{u - d},$$

respectively. (To avoid arbitrage we must have  $u - c > R > d - c$ .) These risk-neutral probabilities should be used to evaluate securities or ventures that are derivative to the commodity.

**Example 15.9 (A foreign currency put)** Mr. Smith, a successful but cautious U.S. businessman, has sold a product to a Japanese firm, and he will receive payment of 10 million Japanese yen in 6 months. Currently the exchange rate  $y$  is \$825 per yen. To protect the value of this anticipated payment, Mr. Smith is considering the purchase

of a 6-month put of 1 million yen at a strike price of \$80 per yen. Mr. Smith wants to compute the fair value of such a put to see whether the market price is reasonable.

To make the calculation, Mr. Smith notes that the U.S. dollar interest rate is 5% while the Japanese yen interest rate is 8%. The interest on marks acts like a proportional dividend or, equivalently, a negative holding cost. The volatility of the exchange rate is 3% per month.

To find the value of the put, Mr. Smith sets up a binomial lattice with six monthly periods, with  $u = e^{.03} = 1.03045$  and  $d = 1/u = .97045$ . The risk-neutral probability for an up move is

$$q = \frac{(1 + .05/12) - d - .08/12}{u - d} = .387.$$

Mr. Smith then evaluates the put with the usual backward process. Specifically, he sets up a lattice of yen prices using the  $u$  and  $d$  factors defined by the volatility. He then sets up a corresponding lattice for put prices. The terminal values are found easily, and other values are found by discounted risk-neutral valuation using the risk-neutral probabilities.

## Brownian Motion Form\*

Suppose a commodity—let's take copper—has a price governed by geometric Brownian motion as

$$dS = \mu S dt + \sigma S dz \quad (15.27)$$

where  $z$  is a standard Wiener process. If an investor buys copper and holds it, there is a proportional storage cost that is paid at the rate of  $c$  per unit time. If at any moment  $t$  the investor holds copper with total value  $W(t)$ , the holding cost can be paid at the rate of  $cW(t)dt$  by selling copper at this rate. The process for the value of copper holdings is therefore

$$\begin{aligned} dW &= \mu W dt + \sigma W dz - cW dt \\ &= (\mu - c)W dt + \sigma W dz, \end{aligned} \quad (15.28)$$

where  $W(0) = S(0)$ . Equation (15.28) can now be regarded as that governing the value of a security with the holding costs accounted for. We might term  $W$  the value of *net copper*, since it is the net value after holding costs.

If we consider an investment opportunity that involves copper, such as an option on copper futures or a real option on a project that involves copper as a commodity (such as a copper mining operation or an electrical equipment project), we can value this opportunity by risk-neutral techniques. We change the process for net copper to risk-neutral form since it is net copper that can be used in constructing a replication of other securities. Specifically, in a risk-neutral setting with interest rate  $r$ , net copper is governed by

$$dW = rW dt + \sigma W d\hat{z}, \quad (15.29)$$

where  $\hat{z}$  is a standard Wiener process.

The appropriate transformation embodied in the foregoing is that from (15.28) to (15.29), which boils down to the change  $\mu - c \rightarrow r$ . This is equivalent to  $\mu \rightarrow r + c$ . Hence the original copper price in a risk-neutral world satisfies

$$dS = (r + c)S dt + \sigma S d\hat{z}. \quad (15.30)$$

This is the equation that should be used for risk-neutral valuation of copper-related investments.

## 15.12 Martingale Pricing\*

Consider any security with a continuous-time price process  $S(t)$ . Suppose that the interest rate is  $r$  and the security makes no payments for  $0 \leq t \leq T$ . The theory of risk-neutral pricing states that there is a risk-neutral version of the process on  $[0, T]$  such that

$$S(0) = e^{-rt} \hat{E}[S(t)], \quad (15.31)$$

where  $\hat{E}$  denotes expectation in the risk-neutral world. We can translate this expression to time  $t_1$  to write

$$S(t_1) = e^{-r(t_2-t_1)} \hat{E}_{t_1}[S(t_2)]$$

for any  $t_2 > t_1$ , where  $\hat{E}_{t_1}$  denotes risk-neutral expectation as seen<sup>4</sup> at time  $t_1$ . We can then rearrange this expression to

$$e^{-rt_1} S(t_1) = e^{-rt_2} \hat{E}_{t_1}[S(t_2)].$$

Equivalently, if for all  $t$  we define

$$\bar{S}(t) = e^{-rt} S(t),$$

then we have the especially simple expression

$$\bar{S}(t_1) = \hat{E}_{t_1}[\bar{S}(t_2)] \quad (15.32)$$

for all  $t_2 > t_1$ .

In general, a process  $x(t)$  that satisfies  $x(t_1) = \hat{E}_{t_1}[x(t_2)]$  for all  $t_2 > t_1$  is called a **martingale**. The expected future value of a martingale is equal to the current value of the process—there is no systematic drift.

Equation (15.32) states that the security price  $S(t)$  deflated by the discount factor from 0 to  $t$  is a martingale under the risk-neutral probability structure.

Furthermore, our results on risk-neutral evaluation imply, in the same way, that the price process  $P$  of any security which is derivative to  $S$  (and which does not generate intermediate cash flows) must also be a martingale under the same probability structure; that is,

$$\bar{P}(t_1) = \hat{E}_{t_1}[\bar{P}(t_2)]. \quad (15.33)$$

---

<sup>4</sup> In equation (15.31) we could write  $\hat{E}_0$ , but the time reference is understood.

This is just a restatement of the risk-neutral pricing formula because we can unscramble (15.33) to produce

$$P(t_1) = e^{-r(t_2-t_1)} \hat{E}_{t_1} [P(t_2)]. \quad (15.34)$$

**Example 15.10 (Forward value)** Consider a forward contract on a security with price process  $S$ . The contract is written at  $t = 0$  with forward price  $F_0$  for delivery at time  $T$ . The initial value of this contract is  $f_0 = 0$ . At time  $t > 0$ , new contracts have forward price  $F_t$ . What is the value  $f_t$  of the original forward contract at  $t$ ?

The function  $f_t$  is a derivative of the security  $S$ ; hence its deflated price must be a martingale in the risk-neutral world. Hence,

$$\bar{f}_t = \hat{E}_t(\bar{f}_T).$$

Equivalently,

$$e^{-rt} f_t = e^{-rT} \hat{E}_t(f_T) = e^{-rT} \hat{E}_t(S_T - F_0). \quad (15.35)$$

The same argument applied to a contract written at  $t$  with forward price  $F_t$  (and value zero) gives

$$0 = e^{-rT} \hat{E}_t(S_T - F_t)$$

or, equivalently,  $\hat{E}(S_T) = F_t$ . Using this in equation (15.35), we find the desired result,

$$f_t = e^{-r(T-t)}(F_t - F_0),$$

which agrees with the formula derived in Section 12.4 by more elementary (but less general) arguments.

The martingale formulation can be used in the binomial lattice framework as well. The analog of equation (15.33) is

$$\bar{P}_k = \hat{E}_k(\bar{P}_j) \quad (15.36)$$

for  $j > k$ , where

$$\bar{P}_k = \frac{P_k}{(1+r)^k}$$

and  $\hat{E}_k$  denotes expectation at  $k$  with respect to the risk-neutral probabilities. Again  $\bar{P}$  is  $P$  deflated by the discount factor. In the binomial framework (15.36) is usually applied a single step at a time, in which case it is identical, once the interest rate terms are made explicit, to the familiar backward discounted risk-neutral recursive evaluation process.

Because of this association with martingales, the risk-neutral probabilities are often termed **martingale probabilities**. However, in this text we generally prefer risk-neutral terminology to martingale terminology.

## 15.13 Axioms and Black–Scholes

The fundamental pricing principles for derivatives lead to some important pricing results in a continuous-time setting. In particular, this approach quickly leads to the

Black–Scholes equation without the need to produce a replicating portfolio, and it provides new insight into pricing principles in the continuous-time framework.

In our setting,  $x$  is a continuous-time market variable governed by an Ito process. The method recognizes that pricing a derivative of  $x$  is an ongoing process, where at each instant  $t$  we must price the situation that will occur at  $t + dt$ . For this purpose a pricing operator  $\mathbb{P}$  is introduced that gives the price at time  $t$  of a payoff at time  $t + dt$ . This operator can be manipulated according to an associated **operational calculus** based on four axioms:

- 1. Pricing a marketed variable:** If  $x$  is the price of an evolving marketed asset that neither pays dividends nor requires holding costs, then  $\mathbb{P}\{x + dx\} = x$ . (This states that the value at  $t$  of a later market price (at  $t + dt$ ) is the current market price.)
- 2. Pricing a constant:** If  $C$  is a constant, then  $\mathbb{P}\{C\} = C(1 - r dt)$ . (This states, first, that a risk-free asset exists and, second, that its price satisfies Axiom 1; namely, risk-free amounts that grow by  $r dt$  are discounted by  $(1 - r dt)$ .)
- 3.  $\mathbb{P}$  is linear.**
- 4.  $\mathbb{P}\{dt\} = dt$ .** Terms of higher order in  $dt$  are ignored.

It should be clear that the first three of these parallel the three pricing principles for derivatives introduced in Section 12.1.

These axioms easily produce an important result. Using Axioms 1, 2 and 3, we have  $x = \mathbb{P}\{x + dx\} = \mathbb{P}\{x\} + \mathbb{P}\{dx\} = (1 - r dt)x + \mathbb{P}\{dx\}$ , since at  $t$  the quantity  $x$  is a constant. Canceling  $x$  from both sides, we obtain

$$rx dt = \mathbb{P}\{dx\}. \quad (15.37)$$

This is a fundamental pricing equation. Now, if  $V(x, t)$  is a value function in the span of marketed assets  $x$  and  $r$ , there must hold

$$rV(x, t)dt = \mathbb{P}\{dV(x, t)\}. \quad (15.38)$$

This is the general continuous-time pricing equation for derivatives.

Let us use this method to derive quickly the standard Black–Scholes equation. As usual, we assume that the price of the underlying security  $x$  is governed by the geometric Brownian motion process  $dx = \mu x dt + \sigma x dz$ . In addition, there is a value function  $V$  with specified terminal value  $V(x, T)$ . From Ito's lemma and  $(dx)^2 = \sigma^2 x^2 dt$ , we have

$$\begin{aligned} dV &= V_t dt + V_x dx + \frac{1}{2} V_{xx}(dx)^2 \\ &= V_t dt + V_x dx + \frac{1}{2} V_{xx} \sigma^2 x^2 dt. \end{aligned}$$

Applying the pricing operator and using the fundamental equation (15.37) to value  $\mathbb{P}\{dV\}$ , we have

$$\mathbb{P}\{dV\} = V_t dt + V_x rx dt + \frac{1}{2} V_{xx} \sigma^2 x^2 dt. \quad (15.39)$$

Substituting this into (15.38) and canceling  $dt$  we obtain the standard Black–Scholes equation:

$$rV(x, t) = V_t(x, t) + V_x rx + \frac{1}{2} V_{xx} \sigma^2 x^2. \quad (15.40)$$

Notice that the quantity  $\mu$  never even enters this derivation.<sup>5</sup>

**Example 15.11 (Intermediate cash flow)** As an example, suppose that a derivative pays cash at a rate of  $h(x, t)$  at time  $t$ . It is easy to see that this should be added to  $dV(S, t)$  so that the result is

$$rV(x, t) = h(x, t) + V_t(x, t) + V_x rx + \frac{1}{2} V_{xx} \sigma^2 x^2. \quad (15.41)$$

## Market Price of Risk

Suppose as usual that  $x$  is governed by the process  $dx = \mu x dt + \sigma x dz$ , where  $z$  is a standard Wiener process. From the pricing equation (15.37) we have

$$rx dt = \mathbb{P}\{dx\} = \mu x dt + \sigma x \mathbb{P}\{dz\}.$$

Thus

$$\mathbb{P}\{dz\} = -\frac{\mu - r}{\sigma} dt = -\lambda dt, \quad (15.42)$$

where the constant  $\lambda = (\mu - r)/\sigma$  is termed the **market price of risk**. Writing it as  $\mu - r = \lambda\sigma$ , we interpret it as the excess return  $\mu - r$  above  $r$  that must compensate for volatility of level  $\sigma$ . Alternatively,  $-\lambda dt$  is the **cost of risk**—the amount by which the effective return  $\mu dt$  is reduced by the presence of  $dz$ . That is,  $\mu - \lambda\sigma = r$ .

The cost  $\mathbb{P}\{dz\}$  will be the same in any derivative driven by  $dz$ . For example, if a derivative is governed by  $dy = \mu_1 y dt + \sigma_1 y dz$ , we infer directly that the market price of risk must be the same as before. Thus

$$\frac{\mu_1 - r}{\sigma_1} = \frac{\mu - r}{\sigma}.$$

**Example 15.12 (Black–Scholes again)** One can quickly derive the Black–Scholes equation: Start with the Ito process for  $f$  from (15.10); use the basic result  $\mathbb{P}\{df\} = rf dt$ , use the value for  $\mathbb{P}\{dz\}$ , and cancel  $dt$  from both sides. [You may wish to try it.]

## 15.14 Summary

The Black–Scholes equation is a partial differential equation that must be satisfied by any function  $f(S, t)$  that is derivative to the underlying security with price process

$$dS = \mu S dt + \sigma S dz,$$

---

<sup>5</sup> A rigorous justification of this process requires deeper analysis. In general an operational calculus gives correct answers under a loosely defined set of assumptions. Further use and generalization of this calculus is found in Section 19.9 of Chapter 19.

where  $z$  is a standardized Wiener process. In particular, the functions  $S$  and  $e^{rt}$  both satisfy the Black–Scholes equation. The price functions of other derivative securities, such as options, satisfy the same equation, but with different boundary conditions.

It is usually difficult or impossible to solve the Black–Scholes equation explicitly for a given set of boundary conditions. It can be solved for the special case of a call option on a stock that does not pay dividends during the life of the option. The resulting solution formula  $C(S, t)$  is called the Black–Scholes formula for the price of a call option. This formula is expressed in terms of the function  $N$ , the cumulative distribution of a standard normal random variable. The function  $N$  cannot be evaluated in closed form, but accurate approximations are available.

The Black–Scholes equation can be regarded as an instance of risk-neutral pricing. Indeed, the value of a derivative security with payoff  $V(T)$  at  $T$  and no other payments can be written as  $V = e^{-rT}\hat{E}[V(T)]$ , where  $\hat{E}$  denotes expectation with respect to the risk-neutral process  $dS = rSdt + \sigma Sd\hat{z}$ .

Delta is defined as  $\Delta = \partial f / \partial S$ . Delta therefore measures the sensitivity of a derivative asset to the changes in the underlying stock price  $S$ . A portfolio can be hedged by constructing it so that its net delta is zero. Delta can also be used to construct a derivative security synthetically, by replication. To do this, one constructs a special portfolio containing the underlying security in sufficient amounts so that its value is equal to the value of delta times the price of the underlying security. The portfolio also contains the risk-free asset (either short or long) in an amount to make the entire portfolio have value equal to the theoretical value of the derivative. The portfolio is rebalanced periodically so that the value continues to track the theoretical value of the derivative closely. Portfolio insurance is an extension of this idea, but it constructs the replicating portfolio with futures contracts on the underlying security rather than with the underlying security itself.

Associated with an option price is an implied volatility that is inferred using the Black–Scholes formula. If all prices followed the Black–Scholes formula, the implied volatilities for various strike prices would all be equal. Instead, it is commonly observed that a graph of implied volatility versus strike price forms the shape of a smile or a skew. This shape reflects the fact that the underlying risk-neutral probability density is not lognormal. In fact, the detailed shape of the curve can be used to construct this risk-neutral probability density. This in turn provides a way to price complex derivatives on that instrument.

There are several ways to compute the value of options or other derivative securities numerically. Monte Carlo simulation is a simple method that is well suited to European-style options, even those that are path dependent in the sense that the final payoff depends on the particular price path of the underlying security as well as the final price itself (as, for example, a call with strike price equal to the average price of the underlying security during the life of the option). A disadvantage of Monte Carlo is that it may require a very large number of simulation runs.

Finite-difference methods approximate the Black–Scholes equation by a set of algebraic equations, which can be solved numerically. The method can treat American- as well as European-style options, but it cannot treat path-dependent options, except in special cases.

Lattice and tree methods are very popular. A disadvantage is that the size of the lattice or tree often becomes very great. Path-dependent options require trees rather than lattices, and hence the number of nodes can become truly enormous.

Many variations of the option concept exist. Formulas for the theoretical prices of some of these exotic options have been devised, but in most cases the prices must be found numerically.

If storage costs are incurred or dividends are received while holding an asset, those will influence the value of securities derivative to that asset. If the storage costs or dividends are proportional to the asset price, the value of a derivative security can be found by properly adjusting the risk-neutral probabilities or, in the continuous-time case, by adjusting the growth coefficient in the risk-neutral process governing the asset.

If intermediate payments are made or costs incurred while holding a derivative security itself, those additional cash flows can, within the binomial lattice framework, be accounted for at each node during the discounted risk-neutral valuation process, as illustrated in Chapter 14. In the continuous-time framework, additional cash flow rates can be entered as an additional term in the Black–Scholes equation.

The risk-neutral valuation equation can be transformed (easily) to martingale form: the price of a derivative deflated by the discount factor defines a martingale process under the risk-neutral probability structure.

The three pricing principles for derivatives introduced in Chapter 12 can be used in a continuous-time framework by introducing a pricing operator that is governed by four axioms, three of which are direct parallels of the pricing principles. The use of these axioms quickly leads to the Black–Scholes equation. This reflects the fact that the Black–Scholes equation is, at root, simply an expression of the basic pricing principles. A useful concept is the market price of risk  $(\mu - r)/\sigma$ . It must be the same for any derivative that is driven by the same Wiener process.

## Exercises

1. (Numerical evaluation of normal distribution  $\oplus$ ) The cumulative normal distribution can be approximated (to within about six decimal places) by the modified polynomial relation

$$N(x) = \begin{cases} 1 - N'(x)(a_1 k + a_2 k^2 + a_3 k^3 + a_4 k^4 + a_5 k^5) & \text{for } x \geq 0 \\ 1 - N(-x) & \text{for } x < 0 \end{cases},$$

where

$$N'(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

$$k = \frac{1}{1 + \gamma x}$$

$$\gamma = .2316419$$

$$a_1 = .319381530$$

$$a_2 = -.35653782$$

$$a_3 = 1.781477937$$

$$a_4 = -1.821255978$$

$$a_5 = 1.330274429.$$

Use this formula to find the value of a call option with parameters  $T = .5$ ,  $\sigma = .25$ ,  $r = .08$ ,  $K = 35$ , and  $S_0 = \$34$ .

- 2. (Perpetual put  $\diamond$ )** Consider a perpetual American put option (with  $T = \infty$ ). For small stock prices it will be advantageous to exercise the put. Let  $G$  be the largest such stock price. The time-independent Black–Scholes equation becomes

$$\frac{1}{2}\sigma^2 S^2 P''(S) + rSP'(S) - rP(S) = 0$$

for  $G \leq S \leq \infty$ . The appropriate boundary conditions are  $P(\infty) = 0$  and  $P(G) = K - G$ .  $G$  should be chosen to maximize the value of the option.

- (a) Show that  $P(S)$  has the form

$$P(S) = a_1 S + a_2 S^{-\gamma},$$

where  $\gamma = 2r/\sigma^2$ .

- (b) Use the two boundary conditions to show that

$$P(S) = (K - G)(S/G)^{-\gamma}.$$

- (c) Finally, choose  $G$  to maximize  $P(S)$  to conclude that

$$P(S) = \frac{K}{1+\gamma} \left[ \frac{(1+\gamma)S}{\gamma K} \right]^{-\gamma}.$$

- 3. (Sigma estimation  $\oplus$ )** Traders in major financial institutions use the Black–Scholes formula in a backward fashion to infer other traders' estimates of  $\sigma$  from option prices. In fact, traders frequently quote sigmas to each other, rather than prices, to arrange trades. Suppose a call option on a stock that pays no dividend for 6 months has a strike price of \$35, a premium of \$2.15, and time to maturity of 7 weeks. The current short-term T-bill rate is 7%, and the price of the underlying stock is \$36.12. What is the implied volatility of the underlying security?

- 4. (Black–Scholes approximation  $\diamond$ )** Note that to first order  $N(d) = \frac{1}{2} + d/\sqrt{2\pi}$ . Use this to derive the value of a call option when the stock price is at the present value of the strike price; that is,  $S = Ke^{-rT}$ . Specifically, show that  $C \approx .4S\sigma\sqrt{T}$ . Also show that  $\Delta \approx \frac{1}{2} + .2\sigma\sqrt{T}$ . Use these approximations to estimate the value of the call option of Example 15.2.

- 5. (Delta)** Using the same parameters as in Example 15.2, find the value of the 5-month call if the initial value of the stock is \$63. Hence estimate the quantity  $\Delta = \Delta C/\Delta S$ . Estimate  $\Theta = \Delta C/\Delta t$ .

- 6. (A special identity)** Gavin Jones believes that for a derivative security with price  $P(S)$ , the values of  $\Delta$ ,  $\Gamma$ , and  $\Theta$  are related. Show that in fact

$$\Theta + rS\Delta + \frac{1}{2}\sigma^2 S^2 \Gamma = rP.$$

7. (Gamma and theta  $\diamond$ ) Show that for a European call or put on a non-dividend-paying stock

$$\Gamma = \frac{N'(d_1)}{S\sigma\sqrt{T}},$$

$$\Theta = -\frac{SN'(d_1)\sigma}{2\sqrt{T}} - rKe^{-rT}N(d_2).$$

[Hint: Use Exercise 6.]

8. (Power of S) Assume we have a non-dividend-paying stock governed by the stochastic differential equation  $dS_t = \mu S_t dt + \sigma S_t dz_t$  and a risk-free asset carrying a fixed interest rate  $r$ . The value  $B$  of the risk-free asset satisfies  $dB = rB dt$ . Consider a European contract that pays the  $k$ th power (where  $k$  is a positive integer) of the stock at time  $T$ . It can be shown that the price of the contract at time  $t$  ( $t \leq T$ ) has the form

$$P(t, S_t) = S_t^k A(T) e^{-\alpha t}.$$

Find  $A(T)$  and  $\alpha$  explicitly.

9. (Mean reversion) A marketed asset's price  $x$  is governed by the mean reverting process

$$dx = \eta(\theta - x(t)) dt + \sigma dz$$

where  $\eta$ ,  $\theta$ , and  $\sigma$  are positive constants and  $z$  is a Wiener process.

- (a) Let  $V(x, t)$  be a given function. Find the Ito process that governs  $V$ .  
 (b) Suppose also that there is a marketed bond satisfying  $dB(t) = rB(t) dt$  with  $r > 0$ . Let  $V(x, t)$  be the value of a derivative of  $x$  and  $B$ . Using the pricing axioms, find the partial differential equation governing  $V$ .

10. (Vega) For a derivative of an asset that follows standard geometric Brownian motion, it may be useful to find the sensitivity of the derivative price with respect to a parameter of the underlying process. In particular, for the parameter  $\sigma$  the corresponding sensitivity is termed **vega** and is

$$\mathcal{V} = \frac{\partial f(S, t)}{\partial \sigma}.$$

For a call option on stock that follows geometric Brownian motion there holds  $\mathcal{V} = S\sqrt{T}N'(d_1)$ , where  $d_1$  is given by equation (15.16b). What is vega for the corresponding put option?

11. (Rho) As with Exercise 10, the sensitivity of a derivative asset with respect to the underlying interest rate  $r$  is termed **rho** and defined as

$$\text{rho} = \frac{\partial f(s, t)}{\partial r}.$$

For a call option on stock following the Black-Scholes framework,  $\text{rho} = KTe^{-rT}N(d_2)$ , where  $d_2$  is defined by equation (15.16c). What is the value of rho for the corresponding put? [Express in terms of  $N(-d_2)$ .]

12. (Great Western CD  $\oplus$ ) Great Western Bank has offered a special certificate of deposit (CD) tied to the S&P 500. Funds are deposited into the account at the beginning of a month and are held in the account for 3 years. Interest is credited to the account at the end of each year, and the amount of interest paid is based on the performance of the S&P 500 index

during the previous 12 months. Specifically,<sup>6</sup> at the end of the first year, if the value of the index at the end of  $k$  months is  $S_k$ ,  $k = 0, 1, 2, \dots, 12$ , the average of the 12-month index values is defined as  $A = \frac{1}{12} \sum_{k=1}^{12} S_k$  and the interest paid is

$$I = \max[0, (A - S_0)/S_0]$$

times the initial account balance. Interest in the following years is computed in the same fashion, with new values of account balance and index values. Assuming that monthly changes in the S&P 500 index can be modeled as geometric Brownian motion with  $\sigma = .20$ , what risk-free rate is equivalent to this CD? [Hint: Try a tree. Use 2-month or 3-month intervals.]

- 13.** (The control variate method) Suppose that it is desired to estimate the expected value of a random variable  $x$ . (This random variable might be the discounted terminal value of a call option on a stock that is following a risk-neutral random process; then the expected value is the value of the option.) One way to do the estimation is to generate numerous samples of  $x$ , according to its probability distribution, and then take the average of the results. A difficulty with this method is that it may take a very large number of samples to obtain satisfactory results. The process can be speeded up somewhat by the use of an additional random variable  $y$  called a **control variate**. The control variate must be correlated with  $x$ , and its expected value must be known. For example, if  $x$  is the terminal value of a call with a down and out feature, we might choose  $y$  to be the terminal value of a similar call without the down and out feature. We can determine the value of  $E(y) = \bar{y}$  by direct methods such as the Black–Scholes formula or a binomial lattice. But we do expect that if the stock should happen to end high on a particular simulation trial, the value of both  $x$  and  $y$  will be relatively high as well. Hence the two variables are correlated.

The estimate  $\hat{x}$  of  $E(x)$  is made with the formula

$$\hat{x} = x_{\text{avg}} + a(y_{\text{avg}} - \bar{y}).$$

Sometimes a small value of  $a$  is selected arbitrarily. However, an optimal value of  $a$  can be estimated as well. Find the value of  $a$  that minimizes the variance of  $\hat{x}$ . (The result will depend on certain variances and covariances.)

- 14.** (Control variate application  $\oplus$ ) Use the control variate method of Exercise 13 to determine the value of a 5-month Asian call option on a stock with  $S_0 = \$62$ ,  $\sigma = 20\%$ , and  $r = 10\%$  and a strike price of  $\$60$ .

- (a) As a control variate use the 5-month standard call option treated in Example 14.3.
- (b) Use  $S_{\text{avg}}$  as a control variate and compare with part (a).

- 15.** (Pay-later options) Pay-later options are options for which the buyer is not required to pay the premium up front (i.e., at the time that the contract is entered into). At expiration, the holder of a pay-later option *must* exercise the option if it is in the money, in which case he pays the premium at that time. Otherwise the option is left unexercised and no premium is paid.

The stock of the CCC Corporation is currently valued at  $\$12$  and is assumed to possess all the properties of geometric Brownian motion. It has an expected annual return of  $15\%$ , an annual volatility of  $20\%$ , and the annual risk-free rate is  $10\%$ .

---

<sup>6</sup> There were some minor changes in the actual formula.

- (a) Using a binomial lattice, determine the price of a call option on CCC stock maturing in 10 months' time with a strike price of \$14. (Let the distance between nodes on your tree be 1 month in length.)
- (b) Using a similar methodology, determine the premium for a pay-later call with all the same parameters as the call in part (a).
- (c) Compare your answers to parts (a) and (b). Do the answers differ; if so why, if not why not? Under what conditions would you prefer to hold which option?
- 16.** (California housing put  $\oplus$ ) Suppose you buy a new home and finance 90% of the price with a mortgage from a bank. Suppose that a few years later the value of your home falls below your mortgage balance and you decide to default on your loan. California has antideficiency judgment legislation that states that the bank can only recover the value of the house itself, not the entire mortgage balance.<sup>7</sup>
- Suppose you take out a 15-year mortgage for 90% of the home price, and suppose that the risk-free rate is constant at 10%. Assume also that the house has a net value to you (perhaps in saved rent) of 5% of its market value each year. Housing prices have a volatility of 18% per year. What is the value of this put option for a loan of \$90? What is the fair value for the interest rate on your mortgage? (Use the small  $\Delta t$  approximation.)
- 17.** (Forest value) Solve Exercise 16 in Chapter 14 assuming that the annual storage cost of cut lumber is 5% of its value.
- 18.** (Mr. Smith's put) Find the value of the put for Mr. Smith described in Example 15.9.

## References

The classic paper of Black and Scholes [1] initiated the modern approach to options valuation. Another early significant contributor was Merton, many of whose papers are collected in [2]. Merton examined many important special cases, such as perpetual options. Details of options trading are given in [3]. Portfolio insurance is discussed in [4, 5]. For the pioneering work on the smile, see [6]. Also see [7]. The Monte Carlo technique is a classic method for evaluating expected value. Its application to options valuation is treated in [8, 9]. A textbook treatment of general finite-difference methods is [10]. Application to options valuation is discussed in [11, 12]. For a discussion of exotic options see [13, 14]. The idea of Exercise 4 is in [15].

1. Black, F., and M. Scholes (1973), "The Pricing of Options and Corporate Liabilities," *Journal of Political Economy*, **81**, 637–654.
2. Merton, R. C. (1990), *Continuous-Time Finance*, Blackwell, Cambridge, MA.
3. *Characteristics and Risk of Standardized Options*, American Stock Exchange, New York; Chicago Board Options Exchange, Chicago; New York Stock Exchange, New York; Pacific Stock Exchange, San Francisco; Philadelphia Stock Exchange, Philadelphia, February 1994.
4. Leland, H. E. (1980), "Who Should Buy Portfolio Insurance," *Journal of Finance*, **35**, 581–594.
5. Rubinstein, M., and H. E. Leland (1981), "Replicating Options with Positions in Stock and Cash," *Financial Analysts Journal*, **37**, July/August, 63–72.
6. Breeden, D. T., and R. H. Litzenberger (1978), "Prices of State-Contingent Claims Implicit in Option Prices," *Journal of Business*, **51**, no. 4, 621–651.
7. Boyle, P. P. (1977), "Options: A Monte Carlo Approach," *Journal of Financial Economics*, **4**, 323–338.
8. Hull, J. C., and A. White (1988), "The Use of the Control Variate Technique in Option Pricing," *Journal of Financial and Quantitative Analysis*, **23**, 237–251.

<sup>7</sup> This is, of course, a simplification of the law.

9. Mitchell, A., and D. Griffiths (1980), *The Finite Difference Method in Partial Differential Equations*, Wiley, New York.
10. Brennan, M., and E. S. Schwartz (1977), "The Valuation of American Put Options," *Journal of Finance*, **32**, 449–462.
11. Courtadon, G. (1982), "A More Accurate Finite Difference Approximation for the Valuation of Options," *Journal of Financial and Quantitative Analysis*, **17**, 697–705.
12. Rubinstein, M. (1991), "Pay Now, Choose Later," *Risk* (February). (Also see similar articles on other exotic options by the same author in subsequent issues of *Risk*.)
13. Hull, J. C. (2008), *Options, Futures, and Other Derivative Securities*, 7th ed., Prentice Hall, Englewood Cliffs, NJ.
14. Brenner, M., and M. G. Subrahmanyam (1994), "A Simple Approach to Option Valuation and Hedging in the Black–Scholes Model," *Financial Analysts Journal*, March/April, 25–28.

# INTEREST RATE DERIVATIVES

**S**ecurities with payoffs that depend on interest rates are called **interest rate derivatives**. Such securities are extremely important because almost every financial transaction entails exposure to interest rate risk—and interest rate derivatives provide the means for controlling that risk. In addition, as with other derivative securities, interest rate derivatives may also be used creatively to enhance the performance of investment portfolios.

Some examples of interest rate derivatives are listed in the next section. These examples illustrate the complexity of the interest rate environment and the range of financial instruments designed to harness that complexity.

The complexity of the interest rate market is reflected in the theoretical structure used for its analysis. Even in the deterministic case, we found that it is necessary to define an entire term structure of interest rates in order to explain bond prices. When uncertainty is introduced, it is necessary to define a randomly changing term structure. We will find, however, that the concepts and methods that we have developed in the past few chapters—namely, risk-neutral pricing, binomial lattices, and Ito processes—can be combined with the ideas of term structure very nicely to develop a coherent approach to the pricing of interest rate derivative securities. The reader should therefore find this chapter quite interesting, both because the topic is itself extremely important in the investment world, and because it brings together and expands much of the previous material.

## 16.1 Examples of Interest Rate Derivatives

Interest rate derivative securities are relevant to many forms of investment. Here are some examples.

- 1. Bonds** Bonds themselves can be regarded as being derivative to interest rates, although the dependency is quite direct. In particular, the price of a risk-free zero-coupon bond with maturity in  $N$  years is a direct measure of the  $N$ -year interest rate. Coupon-bearing bonds can be regarded, as always, as combinations of zero-coupon bonds.
- 2. Bond futures** Futures on Treasury bonds, Treasury notes, and other interest rate instruments are traded on exchanges. These were discussed in Chapter 12.
- 3. Bond options** An option can be granted on a bond. An American call option on a 10-year Treasury bond would grant the right to purchase the bond at a fixed (strike) price within a fixed period of time.
- 4. Bond futures options** More common than actual bond options are options on bond futures. Such options are traded on an exchange that deals with futures on Treasury notes and other interest rate futures contracts. Such options specify delivery of the underlying futures contract.
- 5. Embedded bond options** Many bonds are **callable**, which means that the issuer of the bond has the right to repurchase the bond according to certain terms. (Usually a bond is callable only after a specified number of years.) A call provision can be regarded as an option granted to the issuer, the option being embedded within the bond itself. The issuer of such a bond will find it advantageous to exercise the call option if interest rates fall below those of the original issue. Some bonds are **putable**, which means that the owner of the bond can require that the issuer redeem the bond under certain conditions. Such bonds grant an embedded put option to the bond holder.
- 6. Mortgages** Typically, a home mortgage carries with it certain prepayment privileges, allowing the mortgagee to repay the loan anytime. (Often there is a repayment penalty for, perhaps, the first 2 years.) The repayment privilege is analogous to a call provision in a bond, with the homeowner taking the role of the issuer. Some mortgages have special features such as rates that adjust with prevailing interest rates.
- 7. Mortgage-backed securities** Mortgages are usually packaged together in mortgage pools. A mortgage-backed security is an ownership share of the income generated by such a pool or an obligation secured by such a pool. The individual mortgages in a pool are typically serviced by banks, which receive the monthly mortgage payments and send them to the mortgage owner. For this reason these securities are also termed **pass throughs**. The overall market for mortgage-backed securities is enormous, surpassing that of the corporate bond market.
- 8. Interest rate caps and floors** It is quite common for a financial institution to offer loans to businesses in which the outstanding balance is charged an interest rate that is pegged to a standard, such as the prime rate or the LIBOR<sup>1</sup> rate. Suppose that the loan has a notional value  $A$  and maturity  $T$  and that (for simplicity)

---

<sup>1</sup> The London Interbank Offered Rate (LIBOR) is the rate used for U.S. dollar borrowing through London intermediaries. There are LIBOR rates for various maturities, such as 1 month, 3 months, 6 months, and so on.

payments are due yearly. The standard rate at year  $k$  is  $r_k$ . At the beginning of each year the standard rate is observed, and it defines the interest that applies to the coming year. That is, the payment at the end of the year is determined by the interest at the beginning of the year. However, the business that took out this loan may worry that the standard rate may increase substantially, say, above some value  $r_C$ . This risk can be eliminated by purchasing a **plain vanilla cap** with maturity  $T$  and that at each time point  $k$  pays  $A \max[r_{k-1} - r_C, 0]$ . This payment will compensate for the extra interest (above  $r_C$ ) of the original loan. Adjustable-rate mortgages often have cap and floor features. The interest rate is updated periodically according to an interest rate index, but the charge cannot exceed a certain specified amount each period and may be limited over the life of the mortgage by an overall cap.

- 9. Swaps** A swap is an agreement between two parties to exchange the cash flows of two interest rate instruments. For example, party A may swap its fixed-income stream with party B's adjustable-rate stream.
- 10. Swaptions** The term is short for *swap option*. A swaption is an option on an interest rate swap. Such options are quite popular among corporations wishing to hedge interest rate risk. (See Exercise 9.) For the student, they represent an excellent example of how the interest rate market is becoming ever more sophisticated.

## 16.2 The Need for a Theory

Wise investors take interest rate movements into account as a form of risk. To analyze this risk systematically, it is best to develop a model of interest rate fluctuations. Development of a model may seem difficult because the interest rate environment is characterized at any one time, not by a single interest rate, but by an entire term structure, composed of a series of spot rates, or a spot rate curve. This entire curve varies its shape with time.

A simplistic approach to modeling the fluctuations is to assume that the individual spot rates move independently of one another in a completely random fashion. This is perhaps acceptable abstractly, but it is not in accord with the observation that rates for adjacent maturities tend to move together. A realistic theory would account for this observation and build additional structure into the model of allowable fluctuations. However, as soon as a specific model is proposed, a new issue arises—that of potential arbitrage.

To see how this issue arises, let us hypothesize, as a simple model restricting the fluctuations, that the term structure is always flat, but that it moves randomly up and down—all rates moving together by the same amount. This simple model was in fact used in the immunization analysis of Chapter 3. To complete the model we could decide on a probabilistic structure for the up and down movements, assuming either a discrete set of possible jumps or a continuous distribution of movements. For the present argument, however, we do not need to be that specific. No matter how the probabilities are assigned, this simple model of term structure variations implies that arbitrage opportunities exist. The simplest proof of this is to look

again at Chapter 3, Example 3.10, which treats the immunization problem of the X Corporation. According to that example, if interest rates are flat at 9%, one can form a portfolio by buying \$292,788 worth of bond 1 and \$121,854 worth of bond 2 while shorting \$414,642 worth of a zero-coupon bond that matures in 10 years. The total cost of this portfolio is zero. However, if the term structure moves either up or down, the net value of the portfolio will increase. Hence there is a chance that a positive profit can be made from the portfolio and no chance of a loss—a classic type B arbitrage situation. (This is a general result for the flat term structure assumption, as shown in Chapter 3, Exercise 16.) This example shows that one cannot arbitrarily select a framework for term structure fluctuation if arbitrage opportunities are to be avoided. How can we find a realistic framework that is arbitrage free?

## 16.3 The Binomial Approach

Our familiar tool—the binomial lattice—provides a suitable framework for constructing interest rate models. We set up a lattice with a basic time span between successive nodes equal to the period we wish to use for representing the term structure—perhaps a week, a month, a quarter, or a year. We then assign a **short rate** (that is, a one-period rate) to each node of the lattice. The interpretation of this lattice is that if the process reaches a specific node, then the one-period rate, for the next period, is the rate specified at that node. To complete the model we may assign probabilities to the various node transitions so that we have a full probabilistic process for the short rate. However, *real* probabilities for node transitions are not relevant for the pricing theory that follows. Instead we will also *assign* a set of risk-neutral node transition probabilities. The assignment of the short rate values and the corresponding risk-neutral probabilities completely defines an interest rate structure for all maturities, as will be demonstrated shortly. It is important to understand that the risk-neutral probabilities are assigned in this case rather than derived from a replication argument.

Since the risk-neutral probabilities are assigned, rather than computed, it is convenient to set them all equal to one-half. We follow this convention in this section. It is convenient as well to establish an indexing convention for the nodes of the lattice. For this purpose, it is easiest to draw the lattice in the right-triangle form shown in Figure 16.1. Note that at time  $t$  there are a total of  $t + 1$  nodes, indexed by  $i$  from 0

**FIGURE 16.1 Indexing system for short rate lattice.** Nodes are double indexed in the form  $(t, i)$ . The  $t$  refers to time as shown at the bottom of the lattice, and  $i$  refers to the height above the lowest part of the lattice.



to  $t$ . A convenient way to visualize this notation is to imagine that the two branches leading from any node are considered to be “up” and “flat.” The index  $i$  at time  $t$  denotes how many ups it has taken to reach the node. A specific node in the lattice is indexed by the pair  $(t, i)$ , with  $t$  being time and  $i$  being the node index at that time. At a node  $(t, i)$  there is specified a short rate  $r_{ti} \geq 0$ , which is the one-period rate at that point.

This lattice forms the basis for pricing interest rate securities by using risk-neutral pricing. When the process is at any node, the value of any interest rate security depends only on that node; and we assume that all node values are related by the risk-neutral pricing formula. For example, consider a given node  $(t, i)$  somewhere in the middle of the lattice, and any interest rate security. Suppose the value of this security at node  $(t, i)$  is  $V_{ti}$ . Then according to the rules of the lattice, this value is related to the value of the security at the next two possible successor nodes according to the risk-neutral pricing formula

$$V_{ti} = \frac{1}{1 + r_{ti}} \left( \frac{1}{2} V_{t+1,i+1} + \frac{1}{2} V_{t+1,i} \right) + D_{ti}, \quad (16.1)$$

where  $D_{ti}$  is the dividend payment<sup>2</sup> at node  $(t, i)$ .

## Implied Term Structure

It may seem that we are a long way from having specified an entire term structure model, since all we have are short rates—but actually the whole structure is already there. We just have to extract it. The extraction is accomplished in the same way that a spot rate curve is extracted from a series of one-period forward rates in the deterministic case. For the binomial lattice, the extraction is based on risk-neutral pricing. To see how this works, suppose that we are at the initial time, at node  $(0, 0)$ . The one-period spot rate is simply  $r_{00}$ , as defined at that node. To find the two-period spot rate, we consider a bond that pays \$1 at time 2. We find its value in two steps, working backward using the risk-neutral pricing formula. In detail, suppose for simplicity the period length is a full year. Denote the price at node  $(t, i)$  of the bond that matures at year 2 by  $P_{ti}(2)$ . Then,

$$\begin{aligned} P_{10}(2) &= \frac{1}{1 + r_{10}} \left( \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \right) = \frac{1}{1 + r_{10}} \\ P_{11}(2) &= \frac{1}{1 + r_{11}} \left( \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \right) = \frac{1}{1 + r_{11}} \end{aligned}$$

and, next,

$$P_{00}(2) = \frac{1}{1 + r_{00}} \left[ \frac{1}{2} P_{10}(2) + \frac{1}{2} P_{11}(2) \right].$$

---

<sup>2</sup> This formula assumes that  $D_{ti}$  depends only on  $t$  and  $i$ . For some complex securities, this does not hold and the valuation process is then path dependent. Such cases are illustrated in later sections.

**FIGURE 16.2 Simple short rate lattice and valuation of a 4-year bond.** The bond is valued by working backward in the lower lattice, starting from the terminal value of 1.0 and discounting with the short rate values in the upper lattice.

		Short rate		
		.2600	.2000	.1800
			.1538	.1384
			.1183	.1065
			.0910	.0819
			.0700	.0630
		Bond value		
			1.0000	
			.8667	1.0000
			.7916	.9038
			.7515	.8481
			.7334	.8180
			.8909	.9514
				1.0000

This process can be applied to evaluate the price  $P_{00}(k)$  for any  $k$ . The corresponding spot rate for period  $k$  is then the rate  $s_k$  that satisfies

$$\frac{1}{(1+s_k)^k} = P_{00}(k).$$

**Example 16.1 (A simple short rate lattice)** Figure 16.2 shows a short rate lattice giving the rates for 6 years. (The period length is 1 year.) The figure was constructed by using an up factor of  $u = 1.3$  and a flat (or down) factor of  $d = .9$ . Risk-neutral probabilities for the lattice were assigned as  $q = .5$  for up and  $1 - q = .5$  for flat.

The entire term structure of interest rates can be determined from this lattice by computing the prices of the zero-coupon bonds of various maturities. An example of such a calculation is shown in the lower part of the figure for a bond maturing at time 4. The value is computed by moving backward through the lattice in the familiar way, at each period weighting the next period's values by the risk-neutral probabilities and discounting by the one-period rate. For example, the top entry in the third column is  $P_{22}(4) = \frac{1}{2}(.8667 + .9038)/1.1183 = .7916$ . The value of the bond at time zero is found to be  $.7334$  times its face value. This corresponds to a spot rate from time zero to time 4 of  $s_4 = (1/.7334)^{-25} - 1 = .0806$ . The other spot rates can be calculated in a similar way by constructing a lattice of the corresponding length with 1's in the final column. If this is done, the resulting term structure is found to be  $(.0700, .0734, .0769, .0806, .0844, .0882)$ . Note how the term structure rises smoothly in a manner that is fairly characteristic of actual term structures.

A short rate binomial lattice gives birth to a whole family of spot rate curves, depicting the way the term structure varies randomly with time. To see this, imagine the process initially at the node  $(0,0)$ . The corresponding term structure (spot rate curve) can be determined by the calculations illustrated in the foregoing example. After one period the process moves to one of the two successor nodes. This successor node is then considered to be the new initial node of a (smaller) short rate lattice that is a sublattice of the original one. A corresponding spot rate curve can be computed

exactly as before, but it will have somewhat different values, representing the one-period change. If the process had moved to the other possible node, the corresponding spot rate curve would be somewhat different still. We can therefore visualize a spot rate curve associated with every node in the lattice. As the underlying process moves from node to node, the entire spot rate curve changes.

## No Arbitrage Opportunities

Is the term structure determined from the short rate binomial lattice free from arbitrage possibilities? Yes! This important fact follows from the risk-neutral pricing formula. To prove it, first consider the possibility of arbitrage over a single period, starting at node  $(t, i)$ . Any security at that node is defined by its values  $V_{t+1, i}$  and  $V_{t+1, i+1}$  at period  $t + 1$  and its price  $P_{ti}$  at  $(t, i)$ . These are related by

$$P_{ti} = \frac{1}{2} \frac{V_{t+1, i} + V_{t+1, i+1}}{1 + r_{ti}}.$$

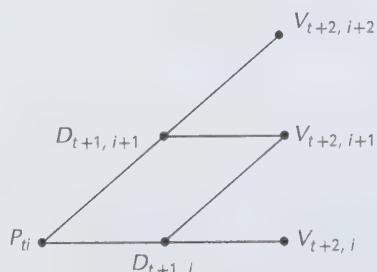
If this security represents an arbitrage, then we must have  $P_{ti} \leq 0$  and  $V_{t+1, i} \geq 0, V_{t+1, i+1} \geq 0$  with one of these inequalities being strict. This is clearly impossible since all coefficients in the equation linking these values are positive. Hence no arbitrage is possible over one period.

The argument for two periods is similar. A security will have price  $P_{ti}$  at time  $t$ , payouts  $D_{t+1, i}, D_{t+1, i+1}$  at time  $t + 1$ , and values  $V_{t+2, i}, V_{t+2, i+1}, V_{t+2, i+2}$  at time  $t + 2$ . It should be clear (see Figure 16.3) that these values are related by

$$P_{ti} = \frac{1}{2} \frac{D_{t+1, i} + D_{t+1, i+1}}{1 + r_{ti}} + \frac{1}{4} \frac{V_{t+2, i} + V_{t+2, i+1}}{(1 + r_{t+1, i})(1 + r_{ti})} + \frac{1}{4} \frac{V_{t+2, i+1} + V_{t+2, i+2}}{(1 + r_{t+1, i+1})(1 + r_{ti})}.$$

Again, for an arbitrage, all variables on the right must be greater than or equal to zero, and  $P_{ti}$  must be less than or equal to zero, with at least one strict inequality. Clearly this is not possible. Hence no two-period arbitrage exists. The argument can be extended to an arbitrary number of periods. Therefore the short rate lattice approach to modeling interest rates is arbitrage free, and hence specification of a short rate lattice provides a workable model of interest rate variations.

**FIGURE 16.3 No arbitrage is possible.** The initial price  $P_{ti}$  is determined by discounted risk-neutral valuation. If all payoffs are non-negative, then the initial price must also be nonnegative.



## 16.4 Pricing Applications

Many interesting securities can be priced with the short rate lattice. Sometimes the short rate lattice together with the promised payout pattern on the nodes of the lattice is all that is needed to set up a backward calculation to determine value. Other times somewhat more subtle techniques must be used. But a wide assortment of problems are amenable to fairly quick calculation using the binomial lattice framework. This section discusses and illustrates a representative group of important and interesting applications of this type.

### Bond Derivatives

The previous section showed how to calculate the value of zero-coupon bonds using the binomial lattice methodology. It is a straightforward extension to calculate the value of other bonds. To calculate the value of a derivative of a bond, we proceed in two steps: first we calculate the price lattice of the bond itself, then we calculate the value of the derivative. We illustrate the procedure for an option on a bond.

**Example 16.2 (A bond option)** Consider a zero-coupon bond that has 4 years remaining to maturity and is selling at a current price of 73.34. Suppose that we are granted a European option to purchase this bond in 2 years at a strike price of 84.00. What is the value of this option?

We assume that the term structure is governed by the short rate lattice of Example 16.1. The value of the zero-coupon bond at any node is indicated in the bond price lattice shown in the bottom portion of Figure 16.2. To evaluate the option we only need the first three periods of this lattice. The value at expiration of the option is  $\max(0, P - K)$ , where  $P$  is the price of the bond at expiration and  $K$  is the strike price. We can then construct a small lattice to determine the option value, as shown in Figure 16.4. The last column shows the value of the option at expiration. The earlier columns show the value obtained by working backward (as usual), using the risk-neutral probabilities of .5 and discounting according to the corresponding values in the short rate lattice. We conclude that the value of the option is 1.4703.

### Forwards and Futures\*

Forward and futures contracts on interest rate securities, such as bonds, are easily treated by the binomial lattice method. This method provides additional insight into the results of Chapter 12 and generalizes those results in important ways, since it is

**FIGURE 16.4 Bond option calculation.** The standard backward method is applied.

			0
	.3712	.81	
1.4703	2.7752	5.09	

not necessary to assume that interest rates are deterministic. Actually, the results for forward contracts are not influenced much by the introduction of uncertainty, but the results for futures are. This means, in particular, that the futures–forwards equivalence result no longer applies. However, the calculations required for interest rate futures are simple.

**Example 16.3 (A bond forward)** Consider a forward contract to purchase a 2-year, 10% Treasury bond 4 years from now. Assume that the short rate process follows the lattice of the previous examples, as shown in Figure 16.2; and assume that coupons are paid yearly and that the contract specifies that delivery will be made just after the coupon payment at the beginning of year 4.

The first step of the calculation is to find the value of the Treasury bond at the beginning of the fourth year. This is done in the usual way by backward calculation, as shown on the right side of Figure 16.5. In the calculation the coupon payments for years 5 and 6 are included. For example, the top entry in year 5 is  $\frac{1}{1.26}(.5 \times 110 + .5 \times 110) + 10 = 97.31$ . The column for year 4 is computed in a similar way, but without the coupon. The figures in the column for year 4 are the prices that the bond would sell for that year.

The left part of the lattice continues the backward calculation, but does not include any coupon payments. The resulting value at the initial node is the value of the 2-year bond delivered at year 4, but paid for at year zero. This is 72.90.

With the forward contract there is no initial payment; the payment is at year 4. This delay of payment has time value, which is determined by the value of a 4-year zero-coupon bond. The value of such a zero-coupon bond was calculated in Example 16.1 to be 73.34. We can find the correct forward price of the bond by comparing it with the forward price of \$100 cash that is to be delivered in 4 years; this forward price is of course just \$100. Hence the correct price of the forward is

$$\begin{aligned} F_0 &= \text{forward price of bond} \\ &= \text{forward price of } \$100 \times \frac{\text{current value of bond}}{\text{current value of } \$100} \\ &= 100 \times \frac{72.90}{73.34} = 99.40. \end{aligned}$$

**FIGURE 16.5 Lattice for bond forward.** The value of the bond is calculated backward from year 6 to year 4. The forward price is then computed backward using the year 4 bond values as final values.

		Year				
0	1	2	3	4	5	6
		Forward period			Bond	
						110
						97.31
					83.56	103.23
					92.69	107.82
				76.38	105.5	111.27
			73.07	87.06	113.80	110
		72.2	84.46	95.69	109.7	115.63
72.9	83.81	93.72	102.4			110

**FIGURE 16.6 Lattice for bond future.** Futures prices are computed by averaging in backward steps without discounting.

Year				
0	1	2	3	4
Futures period				
			88.13	83.56
		92.23	96.33	92.69
	95.88	99.54	102.75	99.69
99.12	102.36	105.18	107.61	105.53
				109.68

## Futures\*

The pricing of futures contracts is also easy using a binomial lattice. The method is best described by a continuation of Example 16.3.

**Example 16.4 (A bond future)** Consider a futures contract on the 2-year, 10% bond to be purchased in 4 years. As before, we need to know the value of the bond at each node for year 4, when the futures contract is due. This calculation was carried out in the previous example, and we simply enter the values in a new lattice at year 4, as shown in Figure 16.6. Now suppose that you are at the top node of year 3, and that the price of the futures contract is  $F$  at that point. You pay nothing then, but next period you would obtain a profit of either  $83.56 - F$  or  $92.69 - F$ . The price you should pay at year 3 is therefore  $.5(83.56 - F) + .5(92.69 - F)$ , discounted by the short rate at that point. But this price is zero, since you pay nothing for the contract. Hence  $F = .5(83.56 + 92.69) = 88.13$ . In other words, the futures price is the average of the two next prices (using the risk-neutral probabilities). This argument can be applied to every previous node. So we just work backward, computing averages *without* discounting. The value at the initial node is the price of the futures contract; namely, 99.12. Note that indeed this value is slightly different than the corresponding forward price of 99.40, thus demonstrating that futures-forward equivalence does not hold when interest rates are random (although the discrepancy is likely to be small).

## 16.5 Leveling and Adjustable-Rate Loans\*

Luckily we have been able to solve most pricing problems in this book using binomial lattices, rather than more complex tree structures. Lattices are very desirable since the number of terminal nodes in a lattice grows only in proportion to  $n$ , the number of periods, whereas for more general trees the number of nodes may grow geometrically (such as  $2^n$  for a binomial tree). Hence if a lattice can be used, representation will be relatively easy and computational effort will be relatively small; whereas everything is more difficult if a full tree is required. Not surprisingly, we are willing to work hard to convert tree structures into lattice structures when that is possible. This section describes a method for doing just that, and then applies the method to the evaluation of adjustable-rate loans.

When using a lattice, nodes are typically defined by the value of some underlying variable that uniquely determines the cash flow at that node. For example, for standard

options, the stock price serves that function, whereas for a bond the short rate is used. If the cash flows associated with a node depend on the path used to arrive at the node, then the cash flow process is said to be **path dependent** and, as discussed in Section 15.10, a lattice is not an appropriate structure. A tree structure, on the other hand, does not have this shortcoming because each node in a tree is reached by a unique path. Hence one way to solve path-dependent problems is to separate all the combined nodes in a lattice, thereby producing a tree that represents the same problem.

Usually, what is going on in a path-dependent case is that more than one variable is needed to describe the cash flow at a node. Sometimes we can collapse these variables into one and salvage the lattice.

We term the technique that we use **leveling** for a reason that will become clear. It applies to situations where cash flow is defined by two variables, say,  $j$  and  $x$ . The first of these is a discrete variable that by itself would define a lattice. The second variable is a continuous variable that is also needed to define cash flow. As an example, consider the Complexico gold mine with random gold prices (which was treated in Chapter 14, Example 14.8). The gold price can be modeled as a binomial lattice, so this price serves as the lattice variable  $j$ . However, after arriving at a lattice node, the cash flow there depends also on the amount of gold remaining in the mine, and hence this amount serves as the  $x$  variable. The mine value is path dependent because the amount  $x$  at any gold price node depends on the path that led to that node. Problems of this type look discouraging because we fear that we might need a lot more nodes to account for the  $x$  dependence.

The path-independent dilemma can be circumvented if the price at a node can be proved to be proportional to the variable  $x$ . If this is the case, we can decide on a fixed level  $x_0$  of  $x$  and use this one level at all nodes, then later scale the results appropriately. Specifically, when working backward, at any node  $j$  we value the security price  $V$  at node  $j$  using the underlying variable values  $j$  and  $x_0$ . The resulting value is  $V_j$ . The step-by-step backward computation is simple because we can easily keep track of the changes in  $x$  for a single step. For example, suppose we are at node  $j$  and we need the price at node  $j+1$ , which is one step ahead, but we need the price at  $j+1$  when  $x \neq x_0$ . By linearity this price is  $(x/x_0)V_{j+1}$ , where  $V_{j+1}$  is the price at  $j+1$  when  $x = x_0$ . Things are especially simple if we choose  $x_0 = 1$ . Then the price at any node  $j$  and level  $x$  is of the form  $V_j(x) = K_j x$ . We just need to keep track of the  $K_j$ 's; then multiply by the appropriate  $x$ .

The method is called leveling because the  $x$  variable is kept at a constant level. The Complexico gold mine problem was solved this way, after it was found that the lease value was linear in the gold reserve amount  $x$ . The method seems to be especially valuable in interest rate derivative problems. We shall use it to treat adjustable-rate loans in the next subsection. That example should clarify the method.

## Adjustable-Rate Loans

Adjustable-rate loans are very common and very important. A typical adjustable-rate loan charges an interest rate in any period that is tied to a standard index, such as the 3-month T-bill rate. For example, the rate charged might be the T-bill rate plus

2 percentage points. However, if the loan is to be amortized over a fixed number of periods (that is, it is to be paid off essentially uniformly), a change in interest rate implies a change in the level of the required payment. The payment in any period is calculated so that the loan will be retired at the maturity date, under the assumption that the interest rate will remain constant until then.

Suppose you were to try to evaluate such a loan. You could take the perspective of the bank that makes the loan, and see how much the bank would pay for the (random) income stream represented by the loan repayment schedule. You would start with a binomial lattice model of the T-bill rate. Then you would be inclined to enter the payments due at any node in the lattice and evaluate this payment structure by backward calculation in the standard way. However, in thinking about this, you would soon discover that the payments could not be entered on the lattice in a unique way because the payment due at any node depends not only on that node, but also on the path taken to get to that node. For example, if a path of high interest rates were taken, the loan balance might be larger than if a path of low interest rates were taken. The loan balance at a node therefore depends on the particular history of interest rates. Your thought at this point would most likely be “Oh, no; it looks like I might have to use a binomial tree, with its thousands of nodes, instead of a lattice. But wait; maybe I can use leveling.”

**Example 16.5 (The auto buyer’s dilemma)** Denise just graduated from college and has agreed to purchase a new automobile. She is now faced with the decision of how to finance the \$10,000 balance she owes after her down payment. She has decided on a 5-year loan, but is given two choices: (A) a fixed-rate loan at 10% interest or (B) an adjustable-rate loan with interest that at any year is 2 points above the 1-year T-bill rate at the beginning of that year. Currently the T-bill rate is 7%. She wants to know which is the better deal.

Denise is pretty adept with spreadsheet programs, so she does a little homework that night. First she decides that the T-bill rate can be modeled by the lattice that we used earlier in Example 16.1. She decides to take the viewpoint of the bank and see what the two loans are worth to it. She makes the assumption that all payments are made annually, starting at the end of the first year.

The fixed-rate loan is easy. The payments are found by using the annuity formula in Chapter 3. Namely,

$$A = \frac{r(1+r)^n P}{(1+r)^n - 1}.$$

For  $P = \$10,000$ ,  $r = 10\%$ , and  $n = 5$  this yields  $A = \$2,638$ , which is the annual payment. The cash flow at each node is shown on the lattice on the left side of Figure 16.7. The lattice on the right side of the figure shows the corresponding value of this cash flow computed using the interest rates of Example 16.1.<sup>3</sup> Denise concludes that the fixed-rate loan is worth \$561.10 to the bank.

---

<sup>3</sup> The loan value can equivalently be calculated as  $-\$10,000 + \sum_{k=1}^5 \left[ \$2,638 / (1 + s_k)^k \right]$ , where the  $s_k$ ’s are the spot rates implied by the short rate lattice.

Year					Year				
0	1	2	3	4	0	1	2	3	4
Payment received					Loan value				
					2,638				2,638
					2,638	2,638			4,836.5
					2,638	2,638	2,638		6,881.3
					2,638	2,638	2,638	2,638	8,914.7
					2,638	2,638	2,638	2,638	11,009.0
					-10,000	2,638	2,638	2,638	9,350.8
						561.1	11,591.7	9,684.8	7,368.0
									5,111.9
									2,638
									5,160.2
									2,638

**FIGURE 16.7 Value of fixed-rate loan.** The lattice on the right is found by standard discounted risk-neutral evaluation using the payments shown in the left lattice.

For the adjustable-rate loan, Denise quickly recognizes that the cash flows are not unique at a node, but depend on the particular path by which the node was reached. She could proceed by constructing a tree and recording at each node both the short rate and the loan balance. Cash flow at the node would be uniquely determined by these two values. Instead, she preserves the lattice structure by using the leveling technique, working with loans of the same balance at every node. She uses a balance value of \$ 100. At each node she calculates the required annual payment to amortize a loan of value \$100 starting at that time and ending at year 5. These values are shown in the lattice on the left side of Figure 16.8. For example, the top element of year 4 shows \$ 122, which is the amount that must be paid at the end of 1 year to clear a loan of \$100 made at an interest rate of  $20.00\% + 2\%$ . Similarly, the initial node shows \$25.71, which is the amount that would have to be paid at the end of each year to amortize a loan of \$100 over the entire 5-year period at a fixed interest rate of  $7\% + 2\%$ . This table is constructed by using the amortization formula. It could be used on an ongoing basis to find the actual payments of the adjustable-rate loan. Denise would simply find the balance of the loan at the node (which depends on the path to the node) and then apply the amount in the lattice as a payment per \$100 of balance. This payment would be made at the end of the then current year.

The lattice on the right side of Figure 16.8 contains at each node the value to the bank of initiating an adjustable-rate loan for \$100 at that node. But the length of the loan is such that it terminates at the end of the original 5-year period. The lattice has the final values of 0 since loans initiated there would be paid back immediately and no interest payments would be received. At the top node of year 4 the bank could loan \$100 at a rate of 22%. This would give it a payment of \$122 next year. This payment has a present value of  $\$122/1.20 = \$101.67$ . Subtracting the \$100 loan outlay gives a net present value profit of \$1.67.

The earlier nodes are a bit more complicated. The top node of year 3 is calculated by noting that a new loan of \$100 will generate a cash flow of \$63.38 next year. Part of this payment is interest payment and part reduces the principal. The remaining principal will be  $\$100 - \$63.38 + (\$15.38 + \$2.00) = \$54.00$ . This principal is received by the bank and then loaned again to Denise during the next period at rates

Year						Year					
0	1	2	3	4	5	0	1	2	3	4	5
Payment						Value per 100					
					100						0
			122	100							1.667 0
		63.38	115.8	100							2.535 1.757 0
	42.95	59.67	111.6	100							3.436 2.665 1.825 0
32.3	40.35	57.13	108.6	100		4.3744	3.601	2.763	1.876		0
25.71	30.39	38.57	55.39	106.6	100	5.349	4.56512	3.723	2.835	1.912	0

**FIGURE 16.8 Value of adjustable-rate loan.** The lattice on the right is found using the leveling technique, keeping the loan balance fixed at \$100. The payments shown in the left lattice are those associated with a balance of \$100.

determined then. (In effect, Denise will pay the bank \$ 63.38 + \$54, and the bank will then issue her a new loan for \$54.) The value of this next loan is either \$1.67 per \$100 or \$1.76 per \$100, each with (risk-neutral) probability of one-half. This amount together with the first payment can be discounted back one period and the \$100 subtracted to obtain the overall net present value of \$2.535. Specifically, the value to the bank of a \$100 loan is

$$\frac{63.38 + 54 + \frac{1}{2}(1.67 + 1.76)54/100}{1.1538} - 100 = 2.535.$$

The first two terms in the numerator are the payments the bank would receive from Denise next period for the loan of \$100. If she really could pay that amount next period, she would have a zero balance then. However, she borrows \$54 next period to make the payment. The value to the bank at that time is captured by the remaining term in the numerator. All this must be discounted at the current rate and the outlay of \$100 subtracted.

Working back through the lattice, Denise finds that a \$100 loan made at year zero is worth \$5.349. Hence the \$10,000 loan is worth \$534.90, which is only slightly lower than the \$561.10 value found for the fixed rate. Hence she concludes that the adjustable-rate loan is somewhat better than the fixed-rate loan in terms of price (although she may wish to carry out a different analysis to see which is best for her utility function, since she is probably unwilling to engage in active T-bill trading to fully hedge the uncertainty).

## 16.6 The Forward Equation

Backward evaluation through a tree or lattice is a powerful method for evaluating financial instruments. There are times when a dual method—a forward recursion—is even better. This forward method is particularly useful for determining the term structure based on a short rate lattice.

In Section 16.4 we saw that a short rate lattice completely determines the term structure. This term structure can be computed by finding the prices of

zero-coupon bonds for each maturity using the backward evaluation method. However, separate recursions and separate price lattices are required for each of these maturities. Hence if there are  $n$  periods,  $n$  separate recursions must be made in order to compute the entire term structure. For large values of  $n$  the number of single-node evaluations is approximately  $n^3/6$ , as compared to  $n^2/2$  for one pass through the entire tree.<sup>4</sup> The forward process described next requires only a single recursion.

The forward recursion is based on calculating **elementary prices**. The elementary price  $P_0(k, s)$  is the price at time zero of a security that pays one unit at time  $k$  and state  $s$ , and pays nothing at any other time or state. The prices  $P_0(k, s)$  are termed elementary prices because they are the prices of elementary securities that have payoff at only one node. We could find  $P_0(k, s)$  for any fixed  $k$  and  $s$  by assigning a 1 at the node  $(k, s)$  in the lattice and then working backward to time zero. Alternatively, we can work forward.

Suppose that elementary prices have been found for all nodes in the lattice for times from 0 through  $k$ . Consider a node of the form  $(k+1, s)$ , where  $s \neq 0, s \neq k+1$ ; that is,  $s$  is not the bottom or the top node of the lattice at time  $k+1$ . This situation is illustrated in Figure 16.9. Such a node has two predecessor nodes (nodes leading to it), namely,  $(k, s-1)$  and  $(k, s)$ . Suppose that a security pays one unit at node  $(k+1, s)$  and nothing elsewhere. If we were to work backward in the lattice, this security would have values  $.5d_{k,s-1}$  and  $.5d_{k,s}$  at the respective predecessor nodes, where  $d_{k,s-1}$  and  $d_{k,s}$  are the one-period discount factors (determined from the short rates at those nodes).

At time zero the values at these two predecessor nodes are worth, by definition of the elementary prices,  $.5d_{k,s-1}P_0(k, s-1)$  and  $.5d_{k,s}P_0(k, s)$ , respectively. The total value at time zero is the sum of these two, and this is the elementary price at  $(k+1, s)$ . Thus  $P_0(k+1, s) = .5d_{k,s-1}P_0(k, s-1) + .5d_{k,s}P_0(k, s)$ . This is a forward recursion because the value at time  $k+1$  is expressed in terms of values at time  $k$ . If  $s = 0$  or  $k+1$ , there is only one predecessor node, and the result is modified accordingly. Overall we obtain the three forms of the forward equation, depending on whether the node is in the middle, at the bottom, or at the top of the lattice,

$$P_0(k+1, s) = \frac{1}{2}[d_{k,s-1}P_0(k, s-1) + d_{k,s}P_0(k, s)], \quad 0 < s < k+1 \quad (16.2a)$$

$$P_0(k+1, 0) = \frac{1}{2}d_{k,0}P_0(k, 0), \quad s = 0 \quad (16.2b)$$

$$P_0(k+1, k+1) = \frac{1}{2}d_{k,k}P_0(k, k), \quad s = k+1.$$

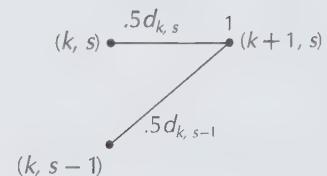
Although we derived this equation through intuitive reasoning, it is possible to derive it algebraically from the backward equation. This forward equation is just a different way of organizing the fundamental risk-neutral pricing equations.

---

<sup>4</sup> A recursion at period  $j-1$  requires  $j$  single evaluations. Hence to evaluate a bond of maturity  $k$  requires  $1+2+\dots+k = (k+1)k/2$  separate evaluations. Since this must be done for all  $n$  maturities, the total is  $\sum_{k=1}^n (k+1)k/2 = [n(n+1)^2/6][1+1/(n+1)]$ . For one pass through the entire tree the number of evaluations is  $n(n+1)/2$ .

**FIGURE 16.9 Construction of forward equation.**

The elementary price for node  $(k+1, s)$  can be expressed as a combination of the elementary prices for the two predecessor nodes.

**FIGURE 16.10 Use of elementary prices to find term structure.**

The elementary prices are determined by a single forward sweep through the lattice. The sum of any column then gives the price of a zero-coupon bond of that maturity. Note that a short rate applies over the coming year while a spot rate and the initial spot rate, although equal, are one column apart.

					.2600
					.2000 .1800
<b>Short rate</b>			.1538	.1384	.1246
		.1183	.1065	.0958	.0862
	.0910	.0819	.0737	.0663	.0597
	.0700	.0630	.0567	.0510	.0459 .0413
					.0069
					.0173 .0468
<b>Elementary prices</b>			.0958	.1754	.2028 .1894
		.2142	.2963	.2757	.2155 .1527
	.4673	.4348	.3046	.1913	.1134 .0648
1.0000	.4673	.2198	.1040	.0495	.0237 .0114
Bond price	.9346	.8679	.8006	.7334	.6670 .6021
Spot rates	.0700	.0734	.0769	.0806	.0844 .0882

The price of any interest rate security can be found easily once the elementary prices are known. We simply multiply the payoff at any node  $(k, s)$  by the price  $P_0(k, s)$  and sum the results over all nodes that have payoffs. For example, the price at time zero of a zero-coupon bond with value 1 that matures at time  $n$  is

$$P_0 = \sum_{s=0}^n P_0(n, s).$$

The forward equation can be used to find the entire term structure corresponding to a short rate tree by a single forward recursion—because all zero-coupon bond prices can be determined.

**Example 16.6 (The simple lattice)** Let us apply the forward equation to Example 16.1. The elementary price lattice can be calculated directly from the short rate lattice. It is shown in Figure 16.10 together with the resulting zero-coupon bond prices and the derived term structure.

As an example of the calculation, both terms in the second column are derived from the single predecessor node; and these terms are equal to one-half times the

discount rate at the first period times the elementary price at 0, which is 1. Hence these values are  $.5/1.07 = .4673$ . The figures directly below the lattice are the sums of the elements above them. These values correspond to prices of zero-coupon bonds. The final figures below the lattice make up the term structure, expressed as spot rates computed directly from the bond prices above. The values agree with those computed in Example 16.1 by the more laborious process.

## 16.7 Matching the Term Structure

Happily we now have an excellent start on a workable methodology for pricing interest rate derivatives, based on the construction of a short rate binomial lattice. From that lattice we can compute the term structure and evaluate interest rate derivatives using the risk-neutral pricing formula and backward recursion. One vital part of this methodology, which we have not yet fully addressed, is how to construct the original short rate lattice so that it is representative of actual interest rate dynamics. This is the subject of this section.

Interest rate fluctuations are similar in character to the fluctuations of stock prices. Therefore a short rate lattice should reflect those basic properties. However, we also know that once a short rate lattice is specified, it implies a certain term structure. It seems appropriate therefore to construct the lattice so that its initial term structure matches the current observed term structure. This is easily accomplished using the concepts and tools developed in the previous sections.

### The Ho–Lee Model

Let us index the nodes of a short rate lattice according to our standard format as  $(k, s)$ , where  $k$  is the time,  $k = 0, 1, \dots, n$ , and  $s$  is the state, with  $s = 0, 1, \dots, k$  at time  $k$ . We must make the assignments  $r_{ks}$  of short rates at each node.

One simple method of assignment is to set

$$r_{ks} = a_k + b_k s. \quad (16.3)$$

This is the Ho–Lee form. It only remains to select the parameters  $a_k$  and  $b_k$  for  $k = 0, 1, \dots, n$ . The variation among nodes at a given time is completely determined by the parameter  $b_k$ . In fact, from any node  $(k - 1, s)$  at time  $k - 1$ , the next rate is either  $a_k + b_k s$  or  $a_k + b_k(s + 1)$ . The difference between the two is  $b_k$ . Indeed, it can be shown easily (see Exercise 6) that the (risk-neutral) standard deviation of the one-period rate is exactly  $b_k/2$ . Hence we refer to  $b_k$  as a **volatility parameter**. The parameter  $a_k$  is a measure of the **aggregate drift** from period 0 to  $k$ . If we remain in state 0, the short rate increases to  $a_k$ .

In the standard Ho–Lee model, the volatility parameters are all set equal to a constant  $b$ , which is characteristic of the observed volatility of interest rates (accounting for the factor of one-half). It therefore remains only to select the  $a_k$ 's; and these can be selected to match the observed term structure at time zero.

If the times are  $0, 1, \dots, n$ , there are  $n+1$  values of  $a_k$  to be chosen and  $n+1$  spot rates to be matched. Hence we have equal numbers of variables and requirements. The only difficulty is that the relation between the  $a_k$ 's and the spot rates is somewhat indirect; but the matching can be carried out numerically.

**Example 16.7 (A 14-year match)** The 12-year term structure used in Example 4.8, Section 4.9, has been extended here to 14 years. We will assume that this is the observed spot rate curve. To match it to a full Ho–Lee model, we must make some assumption concerning volatility. Suppose that we have measured the volatility to be .01 per year, which means that the short rate is likely to fluctuate about 1 percentage point during a year.

We can carry out the match using a spreadsheet package that includes an equation-solving routine. The details are shown in Figure 16.11. The first two lines of the figure show the given spot rates over the 14-year period. The next row shows the parameters  $a_k$  that are used in the Ho–Lee model. These parameters are considered variable by the program. Based on these parameters a short rate lattice is constructed, as shown next in Figure 16.11. From this the forward equations are constructed as another lattice, based on the short rate lattice. The sum of the elements in any column gives the price of a zero-coupon bond with maturity at that date. From these prices, the spot rates can be directly computed. The equation-solving routine is run, adjusting the  $a_k$ 's until the bottom row matches the assumed spot rate values given in the second row.

The spreadsheet method takes advantage of the forward equation and is an appropriate method when the number of periods is not large. When the number of periods is really large, it is better to take advantage of the fact that the spot rate  $s_1$  depends only on  $a_0$ ,  $s_2$  depends only on  $a_0, a_1$ , and so forth. The  $a_i$ 's can therefore be found sequentially by a very rapid process.

## The Black–Derman–Toy Model

An alternative to the model given by equation (16.3) is to assume that the values in the short rate lattice are of the form

$$r_{ks} = a_k e^{b_k s}. \quad (16.4)$$

This can be viewed as a Ho–Lee model applied to  $\ln r_{ks}$ . In this case  $b_k$  represents the volatility of the logarithm of the short rate from time  $k-1$  to  $k$ .

In the simplest version of the Black–Derman–Toy model, the values of  $b_k$  are all equal to a value  $b$ . The  $a_k$ 's are then assigned so that the implied term structure matches the observed forward rates. The computational method is very similar to that for the Ho–Lee model.

## Matching Implied Volatilities

A basic component of the simple Ho–Lee and Black–Derman–Toy models is the volatility parameter. The value of this parameter is usually taken to be the empirical

Year	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Spot	7.67	8.27	8.81	9.31	9.75	10.16	10.52	10.85	11.15	11.42	11.67	11.89	12.09	12.27	
$a$	7.67	8.863	9.878	10.79	11.49	12.18	12.64	13.12	13.5	13.79	14.1	14.23	14.4	14.51	
State	13	12													14.77
	11	10	Short rates												14.75
	9	8													
	7	6													
	5	4													
	3	2													
	1	0													
	14	13	12	Elementary prices											
	9	8	7												
	6	5	4												
	3	2	1												
	0	1	Forward rate												
$\rho_0$	1	.929	.853	.776	.7	.628	.56	.496	.439	.386	.339	.297	.26	.227	.198
	7.67	8.27	8.81	9.31	9.75	10.16	10.52	10.85	11.15	11.42	11.67	11.89	12.09	12.27	

**FIGURE 16.11 Match of term structure.** The observed spot rate curve is given at the top of the figure. Below that are listed some assumed values for the  $a_k$ 's. Using these  $a_k$ 's, the short rate lattice is constructed and the elementary prices are computed by the forward equations. The elementary prices are summed column by column to obtain the zero-coupon bond prices, and these are converted to the forward rates shown in the bottom row. An equation-solving routine is run which adjusts the assumed  $a_k$ 's until the bottom row agrees with the spot rates shown at the top.

value determined from a history of short rate movements. In practice, the resulting lattice often does not accurately predict the observed prices of other interest rate derivatives, such as caps, floors, and swaps. In practice therefore, the volatility is allowed to vary with each time step of the lattice. Thus, as originally suggested, the short rates are of the form

$$r_{k,s} = a_k + b_k s.$$

Each  $b_k/2$  can be considered to be the risk-neutral volatility of the short rate for period  $k$ . These  $b_k$ 's can be determined so that along with the  $a_k$ 's, the resulting risk-neutral lattice will correctly price a set of interest rate caps, floors, or swaps as well as matching the term structure. The volatilities  $b_k/2$  can therefore be regarded as the volatilities implied by the existing market prices of the derivatives.

In particular, for a lattice with  $0 \leq k \leq K$  time periods, the market securities used to determine the implied volatilities might be the  $K - 1$  **plain vanilla caps** of maturities matching the time steps<sup>5</sup>  $k > 0$ . (See Section 16.1.)

## 16.8 Immunization

Our new understanding of interest rate fluctuations and their impact on the term structure provides the basis for a new, more sophisticated approach to bond portfolio immunization, as discussed in Chapters 3, 4, and 5. In those earlier chapters uncertainty was not treated explicitly; instead, a portfolio was immunized against parallel shifts in the spot rate curve. However, we saw in Section 16.2 that the parallel shift assumption is not only simplistic, but in fact inconsistent with a theory that precludes arbitrage. The new approach does not have that weakness.

The new approach is based on the binomial lattice framework. Suppose that we have a series of cash obligations to be paid at specific times in the future, say, up to year  $n$ . Suppose also that we have decided on a specific binomial lattice representation of the short rate. Then we can compute the initial value of the obligation stream using this lattice. One way to compute this value is to first find the term structure at time zero (using the forward equations) and then compute the present value of the obligation stream, just as we learned to do in Chapter 4. Alternatively, but equivalently, we can compute the initial value of the obligation stream by applying the risk-neutral discounting backward process to the obligation stream. The value at the initial node will be the initial (present) value of the stream. To honor the obligation stream, we must have a bond portfolio with this same present value.

After the first period, the value of the obligation stream can take on either of two possible values, corresponding to the values at the two successor nodes. For simplicity assume that no payments must be made at this time. The value at a particular node would correspond to the present value that would be obtained using the new term structure at that node. Likewise, our bond portfolio will have new values at the two successor nodes. Our portfolio is immunized if its value at each of the two successor nodes exactly matches the present value of the obligation at those nodes. In

---

<sup>5</sup> The sequential method for finding  $a_k$  and  $b_k$  can also be used in this case.

**FIGURE 16.12 Initial branching of values.** The initial and next-period values of the two bonds and an obligation are shown. A combination of the bonds will replicate the obligation for one period.

Bond 1		70.96636
	65.95147	71.05353
Bond 2		109.4342
	101.6677	109.497
Obligation		675,949.9
	628,025.6	676,440.4

other words, to immunize for one period, we must match the present values at *three* places—the initial node and the two successor nodes.

The matching might seem complex, but because of the no-arbitrage property of the interest rate structure, things fall into place very nicely. To see how this works, imagine two different bonds that are valued at \$1 at time zero. One of these bonds is the single-period, risk-free bond that pays  $1 + r_{00}$  at each of the two successor nodes. The other is \$1 worth of a zero-coupon bond that matures at year  $n$ . This second bond will have a relatively low value next period if the spot rate increases, but it will have a relatively large value if the spot rate decreases. The two bonds provide two independent outcomes for the next period, and therefore they can be used in combination to replicate the one-period performance of any other interest rate instrument. In particular, they can be combined to replicate the behavior of the obligation.

The solution to the immunization problem is now clear. Using any two dissimilar bonds, we construct a portfolio having the same values at both of the next two states. By the no-arbitrage property, the initial value of this portfolio will be equal to the initial value of the obligation stream that it replicates. Furthermore, the total portfolio consisting of these bonds and the obligation stream is immunized in the sense that its net value is exactly zero initially and at the next period, no matter which state occurs. After one period, the portfolio can be rebalanced to obtain immunization for the next period as well. By continuing to rebalance each period (with the result dependent on the state that occurs), complete immunization over all periods is possible.

**Example 16.8 (Our earlier problem)** We consider again the immunization problem of Example 4.8 in Chapter 4. In this problem we have a \$1 million obligation at the end of 5 years. We wish to immunize this obligation with two bonds. Bond 1 is a 12-year 6% bond with a price of \$65.95. Bond 2 is a 5-year 10% bond with a price of \$101.65. The spot rate curve is known and is equal to that of the Ho–Lee matching problem solved in the last section.

To carry out the immunization we use the short rate lattice found in Example 16.7, since this matches the term structure given in the earlier example. Using this lattice we solve backward for the prices of each of the two bonds and of the obligation. We need to know the results only for the first two periods, which are shown in Figure 16.12. (The initial prices differ slightly from the prices computed earlier due to rounding errors in the lattice.) In each case, the values shown are percentages of the face value.

To construct the immunization, we let  $x_1$  and  $x_2$  be the number of units of bond 1 and bond 2, respectively, in the portfolio. We then solve the equations

$$65.95147x_1 + 101.6677x_2 = 628,025.6 \quad (16.5)$$

$$70.96636x_1 + 109.4342x_2 = 675,949.9. \quad (16.6)$$

(It is not necessary to replicate explicitly state 0 in period 1. This will occur automatically; otherwise there would be an arbitrage opportunity—which is impossible.) The result is that

$$x_1 = 2,165.66 \quad (16.7)$$

$$x_2 = 4,772.38. \quad (16.8)$$

This solution is quite *insensitive* to the volatility assumed when constructing the short rate lattice. Note that the solution is very close to the values of 2,208.17 and 4,744.03 obtained using the standard duration matching method presented in Chapter 4. This seems to be generally true, and hence despite the deeper elegance of the lattice theory, the conventional method of duration matching is frequently used in practice with good results.

## 16.9 Collateralized Mortgage Obligations\*

**Collateralized mortgage obligations** (CMOs) are securities constructed from mortgage pools. The cash flow derived from a pool is sliced up in various ways, and the individual slices define the payout of a particular CMO. The slicing process can be quite intricate, for rather than merely apportioning the principal or the interest payment stream, CMOs are made up of slices that vary the fraction of interest and principal over time. There are numerous variations of the general theme, and new designs are introduced frequently.

The motivating force behind the introduction of CMOs is the prepayment option inherent in real estate mortgages. Homeowners can pay the balance of their mortgage at any time (with some restrictions) and therefore terminate the mortgage. This prepayment feature means that the payment stream of a mortgage is not fixed in advance because the principal might be paid early. This timing uncertainty is somewhat alleviated by the averaging effect derived from a pool, but it is not entirely eliminated because the prepayment pattern cannot be fully predicted. CMOs were devised in order to reduce the variability of the cash flow due to prepayments.

CMOs were first issued by the Federal Home Loan Mortgage Corporation (called Freddie Mac), which buys individual mortgages and forms pools. CMOs issued by Freddie Mac are federally insured against default. Other agencies and corporations now offer CMOs, but those originated by Freddie Mac make up the majority of the market.

The first CMOs were **sequential** CMOs, and they are still very common. In this structure the principal payments are assigned in sequence to different **classes**, or **tranches**, of CMO bonds. Typically there are four to twelve different classes. The total principal of the pool is first divided among the classes. In the early years, mortgage

payments received by the pool are used to pay interest to all classes in proportion to their existing unpaid principal balances, unless they are defined to be Z bonds, in which case owed interest is not paid but instead is accrued and added to the principal balance of that class. The remaining portion of the received mortgage payments is paid to the first class to reduce its principal balance. This continues until the first class is fully retired. After that, the principal of the second class is reduced until it is retired, and so on. Once all previous classes are retired, a Z class bond receives income to reduce its (now greater) principal and to pay interest on that principal.

For example, suppose there are three classes A, B, and Z. Then, as the first mortgage payments are received, interest is paid to classes A and B, and the remaining income is distributed to the A class to reduce its principal. The interest that is due to class Z is paid as principal to class A, thereby speeding the retirement of that class. This foregone interest also augments the principal owed to the Z class. When class A is retired, the principal payments pass to class B, and then finally to class Z. The principal balance patterns are illustrated in Figure 16.13 for a 20-year mortgage pool.

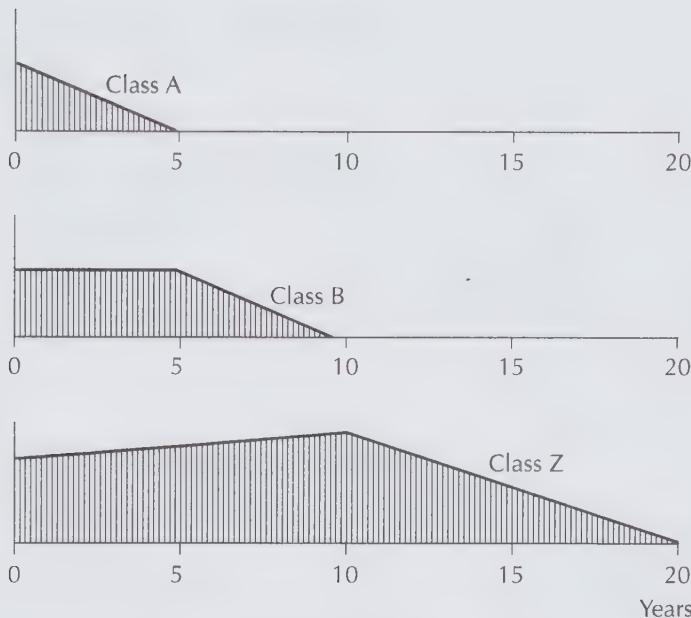
The valuation of CMOs depends very much on the assumed prepayment pattern. A simple approach is to assume a fixed pattern over time. There is in fact a benchmark pattern adopted by the Public Securities Association (PSA). This pattern assumes a prepayment rate of .2% (on an annual basis) the first month, .4% the second month, .6% the third month, and so forth until month 30. After that, the prepayment rate is assumed to be fixed at 6% annually = .5% monthly. For this pattern, or those similar to it, it is easy to project the cash flow pattern for any of the CMO classes. The corresponding value of the CMO class can then be obtained by straightforward discounting using the current spot rate curve. No lattice or tree calculations are required.

In actuality, prepayments depend on prevailing interest rates. Homeowners are more likely to refinance their loans (which entails prepayment of the existing loan) when interest rates are relatively low. Using such a model, a CMO class can be valued using the lattice and tree techniques that we have studied.

**Example 16.9 (Quick, buy this CMO)\*** Mr. Johnathan Quick, the city treasurer of White Falls, is young, well educated, and wants to modernize the financial affairs of the city. A major New York bank has urged him to purchase, for White Falls' account, a portion of class A of a CMO originated by Freddie Mac. This CMO has four classes A, B, C, and Z, each entitled to one-fourth of the principal of a pool of 30-year mortgages carrying an interest rate of 12%. He has been told that these mortgages are guaranteed by the federal government. The current short rate is 10% and the price that he is quoted for the class A bonds is 105.00.

Mr. Quick decides to carry out a simple prototype valuation of this CMO. To do this he first makes a few simple calculations. The yearly payment on a 30-year 12% mortgage is found (see Chapter 3) to be 12.41 per hundred. The interest that will be paid to each of the classes B and C while A is not yet retired is  $25 \times 12\% = 3$ .

He then constructs a short rate lattice covering 4 years, as shown at the top of Figure 16.14. (The lattice starts at the top left node. The successor nodes are the two nodes in the next row.) This lattice has risk-neutral probabilities of .5. Next he assigns estimated prepayment rates. He assigns a 5% annual rate whenever the short rate goes



**FIGURE 16.13 Principal balance patterns of a three-class sequential CMO.** Class A is paid principal before the other classes. When class A is retired, then class B is paid. Class Z does not receive interest until all previous classes are retired. Instead its interest is accrued, augmenting the principal balance.

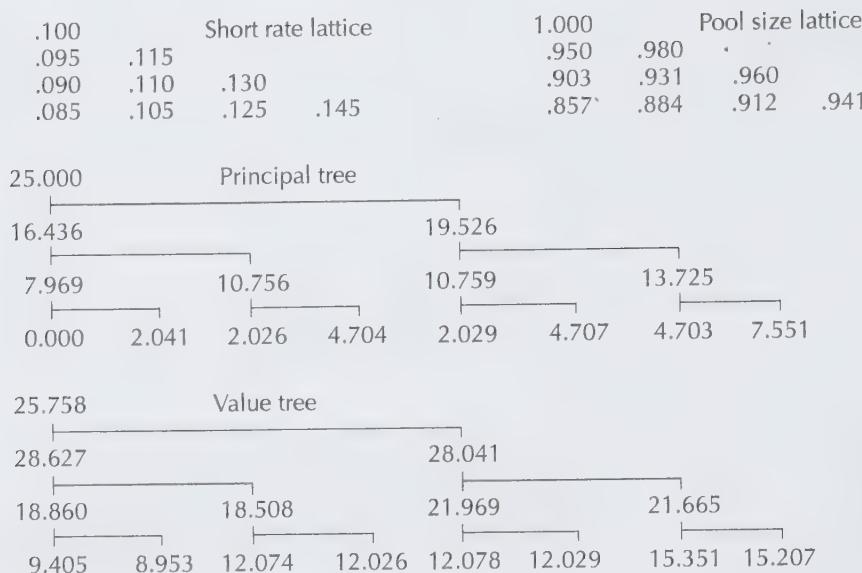
down, and a 2% rate when the short rate goes up. He then puts the remaining pool size fraction on the short rate lattice (shown as a separate array in Figure 16.14).

Quick must keep track of the principal owed to class A. Unfortunately, this principal is path dependent in the original lattice. So he decides that he must use a binomial tree rather than a lattice. He establishes the initial principal to be 25, since class A is entitled to 25% of the total. He arranges his tree in the downward flowing manner, as shown in Figure 16.14. As an example calculation, the final value in the tree is

$$\begin{aligned} & 13.725 \times 1.12 - 12.41 \times .960 + 2 \times 3.00 \\ & - (.960 - .941)[13.725 + 50 + 25(1.12)^3] = 7.551. \end{aligned}$$

In words, the new principal is the old principal times 1 plus the interest rate on the loan (12%), minus the total yearly payment of 12.41 made by the remaining pool, plus the interest payments that must go to classes B and C (but not Z), minus the new prepayment amounts (which is the change in pool size times the total remaining principal). The tree is terminated after 3 years. Mr. Quick assumes that the remaining small amounts of principal will be paid to class A the following year.

To find the value of the class A bond, he uses a tree to carry out backward risk-neutral valuation. A year 3 node value is equal to the year 3 cash flow plus a



**FIGURE 16.14 Quick's CMO valuation.** The top of the figure shows the short rate lattice. Next to that is the lattice showing the corresponding pool size fraction. These lattices start at the top and move downward. A down move is a move directly downward, and an up move is a move downward to the right. Below these is the tree of principal due class A, and finally the corresponding tree of values for class A.

discounted version of next year's principal and interest. The value at an earlier node is equal to the cash flow at that node plus the discounted expected value of the successor node values. For example, the final node value is

$$\begin{aligned} & 7.551 \times 1.12/1.145 + 12.41 \times .960 - 2 \times 3.00 \\ & + (.960 - .941)[13.725 + 50 + 25(1.12)^3] = 15.207. \end{aligned}$$

The final node in the previous row is

$$\begin{aligned} & 12.41 \times .980 - 2 \times 3.00 + (.980 - .960)[19.526 + 50 + 25(1.12)^2] \\ & + .5(15.351 + 15.207)/1.130 = 21.665. \end{aligned}$$

The overall value is 25.758, which when normalized to a base of 100 is  $4 \times 25.758 = 103.032$ . Mr. Quick concludes that the offered price of 105.00 may be a bit high.

He then runs his spreadsheet program again after adding 1 percentage point to each of the short rates and finds the value of 101.112 and therefore concludes that an effective modified duration is  $D_M = 100(103.032 - 101.112)/103.032 = 1.863$  years. This is in accord with the observation that the class A bond is retired very quickly.

Mr. Quick decides to investigate other classes, which he believes may offer substantially greater financial return and whose analyses are sure to offer substantially greater intellectual occupation.

The preceding example shows that the evaluation of CMOs can be quite challenging. If one attempted to carry out the tree methodology of that example, but on a monthly basis and for evaluation of the other classes, very large trees would be required. The main difficulty, of course, is that principal amounts are path dependent. It is for this reason that, in practice, CMO evaluation techniques are usually based on simulation (Monte Carlo) methods. However, it should also be clear from the example that the conceptual principles outlined in the past few chapters are appropriate for this area of finance.

## 16.10 Models of Interest Rate Dynamics\*

In previous sections the short rate was assigned directly by specifying it at every time and state. Although this is a good and practical method, an alternative is to specify the short rate as a process defined by an Ito equation, similar to the processes used to define stock behavior. This allows us to work in continuous time.

In this approach we specify that the (instantaneous) short rate  $r(t)$  satisfies an equation of the Ito type,

$$dr = \mu(r, t)dt + \sigma(r, t)d\hat{z}. \quad (16.9)$$

where  $\hat{z}(t)$  is a standardized Wiener process in the risk-neutral world. Given an initial condition  $r(0)$ , the equation defines a stochastic process  $r(t)$ .

Many such models have been proposed as being good approximations to actual interest rate processes. We list a few of the best-known models:

### 1. Rendleman and Bartter model

$$dr = mr dt + \sigma r d\hat{z}.$$

This model copies the standard geometric Brownian motion model used for stock dynamics. It leads to lognormal distributions of future short rates. It is now, however, rarely advocated as a realistic model of the short rate process.

### 2. Ho–Lee model

$$dr = \theta(t)dt + \sigma d\hat{z}.$$

This is the continuous-time limit of the Ho–Lee model. The function  $\theta(t)$  is chosen so that the resulting forward rate curve matches the current term structure. A potential difficulty with the model is that  $r(t)$  may be negative for some  $t$ .

### 3. Black–Derman–Toy model

$$d \ln r = \theta(t)dt + \sigma d\hat{z}.$$

This is virtually identical to the Ho–Lee model, except that the underlying variable is  $\ln r$  rather than  $r$ . Using Ito's lemma, it can be transformed to the equivalent form

$$dr = [\theta(t) + \frac{1}{2}\sigma^2]rdt + \sigma r d\hat{z}.$$

#### 4. Vasicek model

$$dr = a(b - r)dt + \sigma d\hat{z}.$$

The model has the feature of **mean reversion** in that it tends to be pulled to the value  $b$ . Again, it is possible for  $r(t)$  to be negative, but this is less likely than in other models because of the mean-reversion effect. Indeed, if there were no stochastic term (that is, if  $\sigma = 0$ ), then  $r$  would decrease if it were above  $b$  and it would increase if it were below  $b$ . This feature of mean reversion is considered to be quite important by many researchers and practitioners since it is felt that interest rates have a *natural* home (of about 6%) and that if rates differ widely from this home value, there is a strong tendency to move back to it.

#### 5. Cox, Ingersoll, and Ross model

$$dr = a(b - r)dt + c\sqrt{r} d\hat{z}.$$

In this model not only does the drift have mean reversion, but the stochastic term is multiplied by  $\sqrt{r}$ , implying that the variance of the process increases as the rate  $r$  itself increases.

#### 6. Hull and White model

$$dr = [\theta(t) - ar]dt + \sigma d\hat{z}.$$

This model is essentially the Ho–Lee model with a mean reversion term appended.

#### 7. Black and Karasinski model

$$d \ln r = (\theta - a \ln r)dt + \sigma d\hat{z}.$$

This is the Black–Derman–Toy model with mean reversion.

All of these models are referred to as **single-factor models** because they each depend on a single Wiener process  $\hat{z}$ . There are other models that are **multifactor**, which depend on two or more underlying Wiener processes.

## 16.11 Continuous-Time Solutions\*

The three general methods of solution in discrete time each have a continuous-time analytic counterpart: (1) the method of backward recursion becomes a generalized Black–Scholes partial differential equation, (2) the method of discounted risk-neutral evaluation becomes evaluation of an integral, and (3) the forward recursion method becomes a forward partial differential equation that is dual to the Black–Scholes equation. We shall give some details on the first two of these methods.

## The Backward Equation

The backward equation is perhaps the most useful. Suppose the short rate is governed by the Ito equation (16.9) in a risk-neutral world. And suppose  $f(r, t)$  is a price function for an interest rate security with no payments except at the terminal time. Then it can be shown that  $f$  is governed by the appropriate Black–Scholes equation

$$\frac{\partial f}{\partial t} + \frac{\partial f}{\partial r} \mu(r, t) + \frac{1}{2} \frac{\partial^2 f}{\partial r^2} \sigma(r, t)^2 - rf(r, t) = 0. \quad (16.10)$$

The boundary condition is defined at  $t = T$  and depends on the final payoff structure. This equation is analogous to backward recursion.

For example, suppose we denote by  $P(r, t; T)$  the price at time  $t$  of a zero-coupon bond maturing at time  $T$  when the current short rate (at  $t$ ) is  $r$ . We define the function  $f(r, t) = P(r, t; T)$ , and the appropriate boundary condition is  $f(r, T) = 1$ .

In some cases the backward equation (16.10) can be solved analytically, and this leads to analytic formulas for valuing interest rate derivative securities. In practice, however, numerical solutions are usually required.

**Example 16.10 (Constant interest rate)** The simplest case is when the short rate is governed by  $dr = 0$ , implying that the interest rate is constant. To find the price  $P(r, t; T)$  of a zero-coupon bond, we set  $f(r, t) = P(r, t; T)$ . However, since  $r$  is constant, we may suppress the dependence on  $r$  and write  $f(t)$ . The backward equation reduces to

$$\frac{df(t)}{dt} - rf(t) = 0.$$

This can be written as

$$\frac{df(t)}{f(t)} = r dt$$

or, equivalently, as

$$d \ln f(t) = r dt.$$

This has solution

$$\ln f(t) = c + rt,$$

where  $c$  is a constant. The boundary condition gives  $f(T) = 1$  or, equivalently,  $\ln f(T) = 0$ . Hence we put  $c = -rT$ . The final solution is therefore

$$P(r, t; T) = e^{-r(T-t)}$$

which agrees with what we know about bond values when the interest rate is constant.

**Example 16.11 (A Ho–Lee solution)** As a somewhat more complex example of an analytic solution consider the special case where the short rate is governed by

$$dr = adt + \sigma d\hat{z}.$$

We will try to find the zero-coupon bond price  $P(r, t; T)$ . We set  $f(r, t) = P(r, t; T)$  and solve equation (16.10). Motivated by the solution to the previous example, we try a solution of the form

$$f(r, t) = C(t; T)e^{-r(T-t)}.$$

Substituting this in the Black–Scholes equation, we find

$$\frac{dC(t; T)}{dt} - (T-t) C(t; T)a + \frac{1}{2}(T-t)^2\sigma^2 C(t; T) = 0,$$

where the common factor  $e^{-r(T-t)}$  has been canceled from every term. This leads to the equation

$$d \ln C(t; T) = [(T-t)a - \frac{1}{2}(T-t)^2\sigma^2]dt.$$

Accounting for the boundary condition  $\ln C(T; T) = 0$ , we find

$$\ln C(t; T) = -\frac{1}{2}(T-t)^2a + \frac{1}{6}(T-t)^3\sigma^2.$$

We thus have an explicit formula for  $P(r, t; T)$ .

## Affine Processes\*

A term structure process  $dr = \mu(r, t) dt + \sigma(r, t) d\hat{z}$  is said<sup>6</sup> to be **affine** if both  $\mu(r, t)$  and  $\sigma^2(r, t)$  are affine functions of  $r$ . We shall consider affine processes in the form

$$dr = (\theta - \gamma r) dt + \sqrt{\sigma^2 + \alpha r} d\hat{z},$$

where the  $\theta$ ,  $\gamma$ ,  $\sigma$ , and  $\alpha$  are deterministic functions of  $t$ .

Note that this class includes the term structure models of Ho–Lee; Vasicek; Cox, Ingersoll, and Ross; and Hull and White of Section 16.10. Such a model has a zero coupon price function of the form

$$P(r, t; T) = e^{A(t; T) - B(t; T)r(t)},$$

where  $A(t; T)$  and  $B(t; T)$  satisfy the differential equations

$$\frac{\partial B(t; T)}{\partial t} = B(t; T)\gamma + \frac{1}{2}B(t; T)^2\alpha - 1 \quad (16.11)$$

$$\frac{\partial A(t; T)}{\partial t} = B(t; T)\theta - \frac{1}{2}B(t; T)^2\sigma^2 \quad (16.12)$$

subject to the boundary conditions  $A(T; T) = 0$  and  $B(T; T) = 0$ .

In some cases these equations can be solved in closed form. More generally, a nearly closed-form solution is obtained by the simple process of numerically

---

<sup>6</sup> In general, a function of a variable  $x$  is said to be affine if it is of the form  $a + bx$ .

integrating equation (16.11) backward from  $t = T$ . The resulting solution is then substituted into equation (16.12), which then itself is numerically integrated backward.

**Example 14.12 (Ho–Lee revisited)** For the Ho–Lee model, the equation for  $B(t; T)$  is

$$\frac{\partial B(t; T)}{\partial t} = -1,$$

which has solution  $B(t; T) = (T - t)$ . Using this, the equation for  $A(t; T)$  becomes

$$\frac{\partial A(t; T)}{\partial t} = \theta(T - t) - \frac{1}{2}\sigma^2(T - t)^2,$$

which is easily integrated to obtain

$$A(t; T) = -\frac{1}{2}\theta(T - t)^2 + \frac{1}{6}\sigma^2(T - t)^3.$$

We quickly obtain the complete solution:

$$\ln P(r, t; T) = -\frac{1}{2}\theta(T - t)^2 + \frac{1}{6}\sigma^2(T - t)^3 - (T - t)r.$$

## Risk-Neutral Pricing Formula

The discounted risk-neutral pricing formula also works in the continuous-time case, and it can be used to define the value of any interest rate derivative security. Suppose the security pays a dividend of  $Y(r, t)$  at  $t$ , and suppose that the short rate is governed by the risk-neutral process

$$dr = \mu(r, t)dt + \sigma(r, t) d\hat{z}.$$

Then the value of the security at time zero is

$$v(0) = \hat{E} \left\{ \exp \left[ \int_0^T -r(s) ds \right] Y(r, t) dt \right\} \quad (16.13)$$

where  $\hat{E}$  denotes expectation with respect to the risk-neutral probability defined by the process  $\hat{z}$ . Of course, this formula can rarely be evaluated directly. It does, however, provide a basis for simulation.

## 16.12 Extensions

The short rate process is an elegant, practical, and fairly realistic representation of the stochastic behavior of interest rates. However, it does have limitations because its stochastic nature is represented by a single source of randomness: the Wiener process that drives the short rate. It is a **single-factor model**. Realistically, the interest rate world is more complex than that, and this motivates the use of **multipfactor models**, where the short rate is a function of several factors, each of which is driven by its

own stochastic process. Another modification is the incorporation of processes that have sudden jumps, for such behavior is sometimes observed in the market.

Another approach, proposed by Heath, Jarrow, and Morton, is to formulate the theory in terms of the entire yield curve expressed in terms of the instantaneous forward rates. This HJM model has the advantage of being closer to what traders use, although the analysis and calibration are more complex than for a model based on short rates. A model known as the LIBOR model is similar to the HJM model, except it is formulated in terms of forward rates over finite intervals of time, at points corresponding to the reset times of various instruments, such as interest rate caps and floors.

## 16.13 Summary

Interest rate securities are extremely important because almost every investment entails interest rate risk. Interest rate derivatives, such as bond options, swaps, adjustable-rate mortgages, and mortgage-backed securities, can help control that risk. Analysis of interest rate securities requires a model of term structure variations. Simple models that merely add randomness to a term structure curve are not suitable because they may inadvertently allow arbitrage opportunities.

An elegant and workable approach is to define a short rate lattice spanning several time periods. The rate listed at each node is the interest rate that would apply at that node for loans of one period in length. Two sets of probabilities, are assigned to the arcs of the lattice. The first set defines the *real* probabilities, giving the likelihoods of various transitions. The second set defines the risk-neutral probabilities used for evaluation. Indeed, only the second set is needed for pricing interest rate derivatives. Once the short rate lattice together with the risk-neutral probabilities is constructed, a security such as a bond can be valued by discounted risk-neutral pricing, working backward through the lattice. The short rate at a node defines the discount factor to be used as the process passes through that node.

Seemingly complex securities, such as options on bonds, options on bond futures, and adjustable-rate mortgages, can be evaluated with the discounted risk-neutral approach. In some cases the quantities necessary to determine the cash flow at a node are path dependent, in the sense that these quantities depend on the path to a node as well as on the node itself. In such cases a tree, rather than a lattice, can be used to accurately record the necessary information for the discounted risk-neutral valuation process. However, this can lead to a large increase in the number of nodes. There is a special method termed leveling that transforms an apparently path dependent situation into one that is not path dependent. This method is applicable when the cash flow at a node depends on the node itself and is a linear function of an underlying path-dependent variable. Adjustable-rate loans can be evaluated with this method.

An entire term structure can be extracted from the short rate lattice. One way to do this is to value zero-coupon bonds of all possible maturities. This method requires numerous separate valuation processes. A more efficient way to find the term structure is to construct a lattice of elementary prices. This can be done with a single forward sweep through the original short rate lattice.

The short rate lattice must be constructed carefully in order to give useful results. One common strategy is to construct the lattice so that the term structure that it implies matches the current term structure. Often some volatilities are matched as well. Two of the simplest methods are the Ho–Lee method and the Black–Derman–Toy method.

The short rate lattice also provides a new approach to bond portfolio immunization. In this approach, the portfolio is immunized against initial up and down movements in the short rate.

An important and challenging application of the methodology of interest rate derivative valuation is collateralized mortgage obligations (CMOs). These instruments can have very complex structures, which require careful analysis for proper evaluation. Usually some aspect of their mathematical representation is path dependent, and hence trees or Monte Carlo methods must be employed.

Continuous-time models of the term structure can be constructed by defining a short rate Ito process. This process is driven by a specified risk-neutral standardized Wiener process. Some models of this type lead to analytic expressions for their associated term structure.

## Exercises

- (A callable bond  $\oplus$ ) Construct a short rate lattice for periods (years) 0 through 9 with an initial rate of 6% and with successive rates determined by a multiplicative factor of either  $u = 1.2$  or  $d = .9$ . Assign the risk-neutral probabilities to be .5.
  - Using this lattice, find the value of a 10-year 6% bond.
  - Suppose this bond can be called by the issuing party at any time after 5 years. (When the bond is called, the face value plus the currently due coupon are paid at that time and the bond is canceled.) What is the fair value of this bond?
- (General adjustable formula) Let  $V_{ks}$  be the value of an adjustable-rate loan initiated at period  $k$  and state  $s$  with initial principal of 100. The loan is to be fully paid at period  $n$ . The interest rate charged each period is the short rate of that period plus a premium  $p$ . The loan payment for a period is the amount that would be required to amortize the loan at the charged interest rate equally over the remaining periods. Write an explicit backward recursion formula for  $V_{ks}$  as a function of  $k$  and  $s$ .
- (Bond futures option) Explain how you would find the value of a bond futures option.
- (Adjustable-rate CAP  $\oplus$ ) Suppose that the adjustable-rate auto loan of Example 16.5 is modified by the provision of a CAP that guarantees the borrower that the interest rate to be applied will never exceed 11%. What is the value of this loan to the bank?
- (Forward construction  $\oplus$ ) Use the forward equation to find the spot rate curve for the lattice constructed in Exercise 1.
- (Ho–Lee volatility) Show that for the Ho–Lee model the (risk-neutral) standard deviation of the one-period rate is exactly  $b_k/2$ .
- (Term match  $\oplus$ ) Use the Black–Derman–Toy model with  $b = .01$  to match the term structure of Example 16.7.
- (Swaps) Consider a plain vanilla interest rate swap where party A agrees to make six yearly payments to party B of a fixed rate of interest on a notional principal of \$10 million and in exchange party B will make six yearly payments to party A at the floating short

rate on the same notional principal. Assume that the short rate process is described by the lattice of Example 16.1.

- (a) Set up a lattice that gives the value of the floating rate cash flow stream at every short rate node, and thereby determine the initial value of this stream.
- (b) What fixed rate of interest would equalize both sides of the swap? (Compare with Exercise 11, Chapter 12.)

- 9. (Swaption pricing)** A swaption is an option to enter a swap arrangement in the future. Suppose that company B has a debt of \$10 million financed over 6 years at a fixed rate of interest of 8.64%. Company A offers to sell company B a swaption to swap the fixed rate obligation for a floating rate obligation, with payments equal to the short rate, with the same principal and the same termination date. The swaption can be exercised at the beginning of year 2 (just after the payment for the previous year and when the short rate for the coming year is known). Assuming that the short rate process is that of Example 16.1, how much is this swaption worth?

- 10. (Change of variable  $\diamond$ )** Suppose a short rate process in a risk-neutral world is defined by

$$dr = \mu(r, t)dt + \sigma(r, t)d\hat{z},$$

where  $\hat{z}(t)$  is a standardized Wiener process. A standard way to approximate this equation at a point  $(r, t)$  over a small interval  $\Delta t$  is by the binomial tree shown in Figure 16.15. In this approximation,

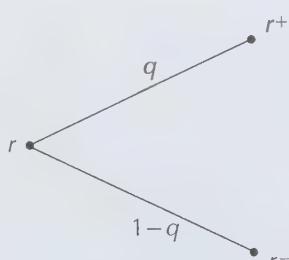
$$\begin{aligned} r^+ &= r + \sigma(r, t)\sqrt{\Delta t} \\ r^- &= r - \sigma(r, t)\sqrt{\Delta t} \\ q &= \frac{1}{2} + \frac{\mu(r, t)\sqrt{\Delta t}}{2\sigma(r, t)}. \end{aligned}$$

- (a) Show that in general this does not produce a recombining lattice. That is, show that an up move followed by a down move is not the same as a down move followed by an up move.
- (b) Consider the change of variable

$$w(r, t) = \int_0^r \frac{dy}{\sigma(y, t)}.$$

Use Ito's lemma to write the process satisfied by  $w(r, t)$ , and show that its volatility term is constant. Conclude that the binomial approximation for  $w(r, t)$  is recombining.

**FIGURE 16.15 Approximation method.** A short rate process can be approximated by a binomial lattice if an appropriate change of variable is used.



- (c) Find the appropriate change of variable for the geometric process

$$dr = \mu r dt + \sigma r d\hat{z}.$$

- 11. (Ho–Lee term structure)** Refer to Example 16.11. Let  $F(t)$  be the forward rate from 0 to  $t$ . By the basic definition of the forward rate, we have the identity

$$e^{-F(t)t} = P(r, 0, t).$$

Find an explicit formula for  $F(t)$ .

- 12. (Continuous zero ⚭)** Gavin wants to dig deep into pricing theory, so he decides to work out an application of Eq. (16.11). He suggests to himself that a simple model of interest rates in the risk-neutral world might be

$$dr = \sigma d\hat{z},$$

where  $\hat{z}$  is standard Brownian motion and where  $r(0) = r_0$ . He is working out a formula for the value of a zero-coupon bond that pays \$1 at time  $T$ , based on Equation (16.11), without using the Black–Scholes equation. Can you? Compare with Example 16.11.

- 13. (Inverse floater)** Consider a bond with a face value of \$1,000 and coupon payment at the end of each period  $k$  given by a rate  $c_k = \max[6\% - r_k, 0]$ , where  $r_k$  is the short rate for period  $k$ . This type of bond is called an **inverse floater**.

The one-period spot rate is currently 4%, and at each period it will either increase 1.5 times or remain constant. The risk-neutral probability implied by the current term structure is 0.5. What is the price of an inverse floater maturing two periods from now? Note there are two coupon payments, with the first being paid one period from now.

- 14. (Forwards and futures)** Consider a payoff  $C$  that will occur in 2 years, taking one of the three possible values  $C_0, C_1, C_2$ . The short rate lattice for these 2 years is shown in Figure 16.16, with  $d_{ij} = \frac{1}{1+r_{ij}}$  being the short rate discount factors and all risk-neutral probabilities being 0.5.

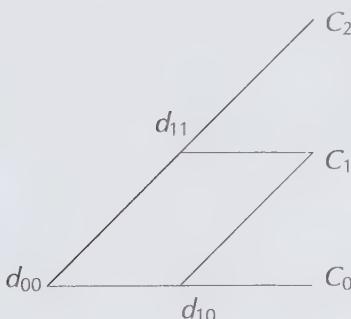


FIGURE 16.16 Short rates.

- (a) What is the futures price  $F_0$  of  $C$ ?  
 (b) What is the forward price  $G_0$  of  $C$ ?  
 (c) Under what conditions will  $F_0 = G_0$  for all possible  $C_0, C_1, C_2$ ?

## References

For general textbook presentations of interest rate derivatives see [1–3]. The forward equation was presented in Jamshidian [4]. The Ho–Lee model was originally developed in [5] without the benefit of the short rate lattice concept. The short rate lattice was used in the presentation of the Black–Derman–Toy model in [6]. For the general theory of affine models see Duffie [1]. The form in this chapter is from [3]. The Heath, Jarrow, Morton model was presented in [7]. For an outline of CMOs and mortgage-backed securities, see [8]. For continuous-time models, see [9–14]. Exercise 10 is based on [15].

1. Duffie, D. (2001), *Dynamic Asset Pricing*, 3rd ed., Princeton University Press, Princeton, NJ.
2. Hull, J. C. (2008), *Options, Futures and Other Derivative Securities*, 2nd ed., Prentice Hall, Englewood Cliffs, NJ.
3. Veronesi, P. (2010), *Fixed Income Securities*, Wiley, New York.
4. Jamshidian, F. (1991), “Forward Induction and Construction of Yield Curve Diffusion Models,” *Journal of Fixed Income*, **1**, 62–74.
5. Ho, T. S. Y., and S.-B. Lee (1986), “Term Structure Movements and Pricing Interest Rate Contingent Claims,” *Journal of Finance*, **41**, 1011–1029.
6. Black, F., E. Derman, and W. Toy (1990), “A One-Factor Model of Interest Rates and Its Application to Treasury Bond Options,” *Financial Analysts Journal*, January–February, **46**, 33–39.
7. Heath, D., R. Jarrow, and A. Morton (1990), “Bond Pricing and the Term Structure of Interest Rates: A Discrete Time Approximation,” *Journal of Financial and Quantitative Analysis*, **25**, 419–440.
8. Fabozzi, F. J., Ed. (1988), *The Handbook of Mortgage-Backed Securities*, rev. ed., Probus, Chicago, IL.
9. Brennan, M., and E. Schwartz (1979), “A Continuous Time Approach to the Pricing of Bonds,” *Journal of Banking and Finance*, **3**, 133–155.
10. Vasicek, O. A. (1977), “An Equilibrium Characterization of the Term Structure,” *Journal of Financial Economics*, **5**, 177–188.
11. Hull, J., and A. White (1990), “Pricing Interest Rate Derivative Securities,” *Review of Financial Studies*, **3**, 573–592.
12. Black, F., and P. Karasinski (1991), “Bond and Option Pricing when Short Rates are Lognormal,” *Financial Analysts Journal*, July–August, **47**, 52–59.
13. Rendleman, R., and B. Bartter (1980), “The Pricing of Options on Debt Securities,” *Journal of Financial and Quantitative Analysis*, **15**, March, 11–24.
14. Cox, J. C., J. E. Ingersoll, and S. A. Ross (1985), “A Theory of the Term Structure of Interest Rates,” *Econometrica*, **53**, 385–407.
15. Nelson, D., and K. Ramaswamy (1989), “Simple Binomial Processes as Diffusion Approximations in Financial Models,” *Review of Financial Studies*, **3**, 393–430.

# 17

## CREDIT RISK

**C**redit risk is the risk associated with nonperformance of the terms of a financial instrument by a counterparty, either through adverse circumstances that lead to default or through willful abrogation. A prime example is the default risk faced by the holders of a firm's debt (in the form of bonds). Another is the risk faced by a bank's depositors when there is a possibility that the bank will fail. Still another is the risk faced by a financial institution when lending money to another company or to a sovereign entity. Adverse credit events include defaults, restructuring, missing of a coupon payment, and any event that causes sudden loss of market value due to the perceived evaluation of credit worthiness. This chapter focuses primarily on default and/or a change in the probability of default. Credit events may be relatively infrequent, but their impact can be huge—even catastrophic. Credit derivatives are designed to mitigate the impact of such events should they occur or, conversely, to bet essentially that such an event does occur. In essence, credit derivatives are similar to insurance policies.

Indeed, some of the most common credit derivatives, called **credit default swaps** (CDS's), parallel almost exactly the structure of an insurance policy. The owner of the policy protects against default risk by making periodic payments (say, quarterly) to the insuring institution and in return receives a payment equal to the loss suffered should default occur within the given time period. This general structure has several variations. However, a major difference in some contracts is that, unlike, say, a fire insurance policy, anyone can purchase a specific traded CDS, even if the purchaser does not own the underlying security. It is like when someone else takes out fire insurance on your house.

The market for credit derivatives has quickly become enormous, with the value of CDS's, for example, being in the tens of trillions of dollars and that of derivatives as a whole being in the hundreds of trillions.

The mathematical study of credit risk and associated credit derivatives represents perhaps some of the most complex analyses in this text. We will introduce new concepts and procedures; yet the methods used in this chapter are, at root, based largely on the fundamental concepts of finance studied in earlier chapters for futures, options, and interest rate derivatives.

There are several quite distinct approaches to the characterization of credit risk. All essentially have the common goal of determining the probability of occurrence of a specified credit event (most often default) or, conversely, the probability of survival. We frequently refer to the survival probability  $p(t)$  as the probability of no default in the time interval  $[0, t]$ . However, this probability must be distinguished as either real or risk neutral; both are used in credit analysis. We use  $p$  for general discussion or for real probability and  $q$  when we mean specifically the risk-neutral probability. It is risk-neutral probabilities that determine prices through risk-neutral expectation and discounting. Real probabilities characterize the risk more directly in the form of variance or measures such as value at risk.

## 17.1 The Classic Merton Model

Modern analytic approaches to credit risk began with the option-based model for a firm's default process proposed by Merton in 1974. This method and extensions of it are termed **structural methods**, for they focus on the structure of a firm's finances and translate that into a measure of risk. Such models continue to provide a foundation for several important credit analysis procedures.

The basic setup is straightforward. Consider a firm that has market value  $V$  consisting of equity and debt, with the debt being in the form of a zero-coupon bond with face value  $K$  and maturity date  $T$ . The market value of the bond at time  $T$  when the firm value is  $V_T$  is

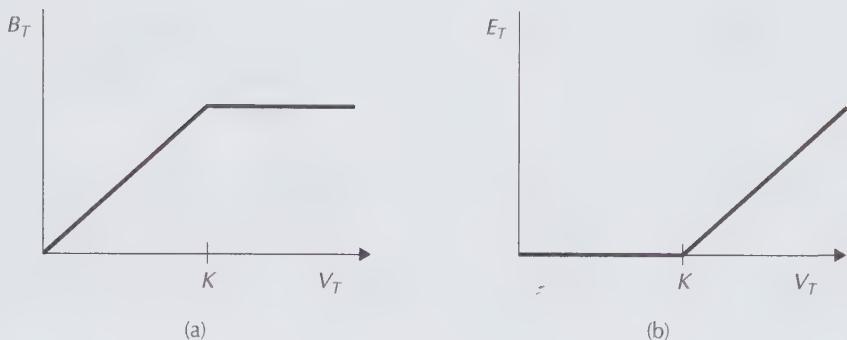
$$B_T = \min(V_T, K) = K - \max(0, K - V_T), \quad (17.1)$$

as shown in Figure 17.1(a). This reflects the fact that if firm value is greater than  $K$ , there is no default and the bondholders receive their promised  $K$ . On the other hand, if at time  $T$  firm value is less than  $K$ , bondholders, having absolute priority of claim, will take over the company and sell it for its value. The formula for bond value can also be considered to be a return of  $K$  and the return of a short put option at strike price  $K$ , as in the second part of equation (17.1).

Similarly, the equity  $E_T$  at  $T$  is

$$E_T = \max(0, V_T - K) = V_T - K + \max(0, K - V_T), \quad (17.2)$$

which in the first part is equivalent to the payoff of a European call option with strike price  $K$ , as shown in Figure 17.1(b). Alternatively, as expressed by the second part



**FIGURE 17.1** (a) Terminal value of a defaultable bond as a function of firm value, and (b) the corresponding value of equity.

of equation (17.2),  $E_T$  can be regarded as the firm having  $V_T - K$  plus a put option with strike price  $K$ . Note that the sum of equity and debt is always equal to  $V_T$ .

If the interest rate  $r$  is constant and the firm's value follows geometric Brownian motion of the form

$$dV = \mu V dt + \sigma V dz \quad (17.3)$$

(where, as usual,  $z$  is a standardized Wiener process), then the value of equity at an intermediate time  $t$  is determined directly by the Black–Scholes formula for a call option from Section 15.3 as

$$E_t = V_t N(d_1) - e^{-r(T-t)} K N(d_2),$$

where

$$d_1 = \frac{\ln(V_t/K) + (r + \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{T-t}}$$

$$d_2 = d_1 - \sigma\sqrt{T-t}.$$

The value of the bond at time 0 can be found by noting that the value at time 0 of the payoff  $K$  is  $e^{-rT} K$  and that the value of the put option can be found from the Black–Scholes formula (determined by put-call parity). Alternatively, the value of the bond can be found from the fact that for all  $t$  there holds  $E_t + B_t = V_t$ . Hence,

$$B_t = V_t - V_t N(d_1) + e^{-r(T-t)} K N(d_2). \quad (17.4)$$

**Example 17.1 (A leveraged firm)** The HiTech firm currently has a value of \$1 million and is financed in part by a 5-year loan that requires payment of \$800,000 at the end of 5 years. The risk-free rate is 5%. If HiTech's value at the end of 5 years is less than \$800,000, then the firm will default and the bondholders will take over the company and sell it for its value. Assume the drift rate of HiTech's value is  $\mu = .15$  and the volatility is  $\sigma = .25$ . The value of the loan can be calculated directly. Using

equation (17.4) with

$$d_1 = [\ln(10/8) + (.05 + .0625/2)5]/(.25\sqrt{5}) = 1.2589$$

$$d_2 = d_1 - .25\sqrt{5} = .566876,$$

we find  $B_0 = \$575,331$  (and, of course, the value actually does not depend on the drift rate  $\mu$ ). This value can be compared to the current value of an \$800,000 5-year zero-coupon bond, which is  $800,000 \exp(-.05 \times 5) = \$623,041$ . The possibility of default reduces the value of the loan.

## Probability of Default

In the context of credit risk, there are other measures, aside from value, that are of interest. One measure is the probability of default, which for the classic model of this section is easily computed. We note first that the geometric Brownian motion model (17.3) for firm value implies (see Section 13.7)

$$\ln V_T = \ln V_0 + \nu T + \sigma z_T,$$

where  $\nu = \mu - \frac{1}{2}\sigma^2$ . Therefore, the probability of default by time  $T$  is

$$\begin{aligned} P[V_T < K] &= P[\ln V_T < \ln K] \\ &= P\left[z_T < \frac{\ln K - \ln V_0 - \nu T}{\sigma}\right] = F_N\left(\frac{\ln(K/V_0) - \nu T}{\sigma\sqrt{T}}\right), \end{aligned} \quad (17.5)$$

where, as usual,  $F_N$  is the distribution function of a standard normal random variable. Indeed, the last equality follows from the fact that  $z_T$  is a normal random variable with variance equal to  $T$ . Formula (17.5) gives the actual probability. The risk-neutral probability is obtained by simply changing  $\mu$  to  $r$  and, hence,  $\nu$  to  $r - \frac{1}{2}\sigma^2$ . (See Section 15.4.)

**Example 17.2 (HiTech continued)** Noticing that the value of the company's bond is somewhat discounted from the current value of a similar zero-coupon value of \$623,041, it seems there is a modest chance that the company will default. We calculate this (real) probability from equation (17.5) (using  $\nu = .15 - \frac{1}{2}(.25^2) = .119$ ) as

$$P[V_T < 800,000] = F_N\left(\frac{\ln .8 - .119 * 5}{.25 * \sqrt{5}}\right) = .072. \quad (17.6)$$

So there is about a 7% chance that the company will default.

## Credit Spread

The **credit spread** is the difference between the promised yield on a defaultable bond and that of an otherwise-equivalent default-free bond. The spread is a convenient and meaningful interpretation of the price difference due to the risk of default.

Suppose at time  $t = 0$  a defaultable zero-coupon bond with maturity date  $T$  has price  $B_T$ . Its yield  $y$  is defined to satisfy the equality  $B_T = \exp(-yT)F$ , where  $F$  is

the face value of the bond. Let  $\bar{B}_T$  and  $\bar{y}$  be, respectively, the price and yield of an equivalent bond that is secure. The spread between the two bonds is defined as

$$S(T) \equiv y - \bar{y} = -\frac{1}{T} \ln(B_T/\bar{B}_T). \quad (17.7)$$

The **term structure of credit spreads** is the curve  $S(T)$  for various lengths of maturity.

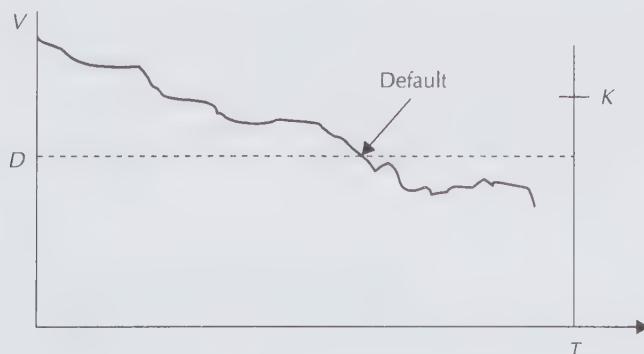
**Example 17.3 (HiTech spread)** For the HiTech case we have  $B = \$575,331$ , and, as found earlier,  $\bar{B} = \$623,041$ . Therefore,

$$S(T) = \frac{1}{5} \ln(623,041/575,331) = .0159,$$

or 159 basis points.<sup>1</sup> The required yield on the defaultable bond is  $5\% + 1.59\% = 6.59\%$ , the increase being necessary to compensate for the risk of default.

## 17.2 First Passage Times

The classic Merton model has several limitations. A major improvement is to recognize that default may occur prior to the maturity date  $T$ . For example, it may be more realistic to define a company to be in default any time its total value falls below the debt level  $K$ . This feature, in turn, can be modified by defining default as occurring whenever firm value falls below a critical value  $D$ , as illustrated in Figure 17.2. The pure Merton model corresponds to setting  $D = 0$  for  $t < T$ ; hence, any value  $D > 0$  will result in a greater default probability than the Merton model. If  $D \geq K$ , bond-holders are fully protected against loss, for in the event of default they take over the



**FIGURE 17.2** When firm value drops below a specified level  $D$ , the company is deemed to be in default.

<sup>1</sup> A basis point is 1/100 of a percent.

company, which is worth at least  $K$ . Normally, however, it is more realistic to use a  $D < K$ , in which case bondholders are exposed to some risk.

The payoff structure, modified this way, is analytically attractive. It can be viewed as containing a down-and-out call option for which a closed form solution is available<sup>2</sup> for  $E_0$ , the value of the equity at  $t = 0$ .

Further modifications of the model have been proposed that may better characterize actual credit events. For instance, the barrier  $D$  may vary with time, the payment to bondholders in the event of default may be subject to delay and various legal and operational costs, there may be several bonds issued by the same company with differing parameter values, and default may entail renegotiation. Finally, the dynamics of firm value may not be adequately characterized by geometric Brownian motion. It is challenging to formulate a model that is an appropriate compromise between analytic tractability and practical realism.

## Lattice Methods

Computational methods, especially simulation and lattice methods, provide potent alternatives to a purely analytic approach. Indeed, many of the features characterizing defaultable bonds can easily be incorporated into evaluation methods that parallel the methods for derivatives developed in previous chapters. We illustrate a few quite simple situations here.

**Example 17.4 (HiTech 2)** Consider the HiTech bond of Example 17.1. The firm value follows geometric Brownian motion with  $V(0) = \$1,000,000$ ,  $\mu = .15$ ,  $\sigma = .25$ ,  $r = .05$ . The bond is described by  $T = 5$ ,  $K = \$800,000$ . According to the method of lattice construction of Section 13.9, the parameters of a corresponding lattice with  $\Delta t = 1$  (not using the small  $\Delta t$  approximation) are  $u = 1.3188$ ,  $d = 1/u$ ,  $p = .7145$ ,  $R = 1.05127$ . Also,  $q = (R - d)/(u - d) = .5227$ . Of course, a lattice with  $\Delta t = 1$  is a rather rough approximation to the continuous-time model, but it illustrates the approach. The value of the bond can be found by straightforward backward evaluation using the risk-neutral probability  $q$  and risk-free return  $R$ .

Figure 17.3 displays the lattices for firm value and for prices of the defaultable bond. The top lattice uses the values of  $u$  and  $d$  to forecast the possible firm values, beginning with \$1,000 (in thousands). The last column of the second lattice shows the

---

<sup>2</sup> For  $D < K$ :

$$E_0 = C(\sigma, T, K, r, V_0) - V_0 \left( \frac{D}{V_0} \right)^{\frac{2r}{\sigma^2} + 1} N(d_1) + Ke^{-rT} \left( \frac{D}{V_0} \right)^{\frac{2r}{\sigma^2}} N(d_2),$$

where

$$d_1 = \frac{(r + \frac{1}{2}\sigma^2)T + \ln(D^2/(KV_0))}{\sigma\sqrt{T}}$$

$$d_2 = \frac{(r - \frac{1}{2}\sigma^2)T + \ln(D^2/(KV_0))}{\sigma\sqrt{T}}$$

<b>Firm value</b>					
1000	1318.80	1739.23	2293.70	3024.93	3989.28
	758.27	1000.00	1318.80	1739.23	2293.70
		574.97	758.27	1000.00	1318.80
			435.98	574.97	758.27
				330.59	435.98
					250.67
<b>Bond Value</b>					
564.34	631.75	684.66	723.87	760.98	800.00
	551.14	641.67	715.27	760.98	800.00
		511.19	630.00	742.04	800.00
			435.98	574.97	758.27
				330.59	435.98
					250.67

**FIGURE 17.3 Firm value and value of a defaultable bond (in thousands of dollars).**  
The standard backward lattice method is used.

final payoff of the bond, which is simply  $\min(V, 800)$ , where  $V$  is the corresponding value of the firm. The initial value of the bond is found to be \$564,340. This compares with \$575,331 found by the Black–Scholes formula for the continuous version. However, a lattice version built with  $\Delta t = 1/4$  gives a bond value of \$572,250, which is quite close to the value determined by the continuous model.

**Example 17.5 (Probability of default)** The probability of default (real or risk neutral) can easily be calculated by constructing a separate lattice. Such a lattice for real probabilities ( $p = .715$ , as found earlier) is shown in Figure 17.4. The final node values are either 0 or 1, depending on whether there is no default or a default, respectively. Earlier node values  $D$  are found by backward recursion of the form  $D = pD^u + (1-p)D^d$  (from the forward nodes), where  $p$  is the actual (real) probability of “up” and  $D^u$  and  $D^d$  are the default probabilities of the two successor nodes. In this case the overall default probability is 14.4%, which is high compared to the 7% given by the continuous-time model. A similar model with  $\Delta t = 1/4$  gives a default probability of 9%.

<b>Default Probability</b>					
0.144	0.073	0.023	0	0	0
	0.322	0.197	0.081	0	0
		0.634	0.489	0.285	0
			1	1	1
				1	1
					1

**FIGURE 17.4 Lattice for default probability.**

## Early Default\*

The lattice method easily accommodates alternative models of default and compensation. For example, Figure 17.5 shows a lattice corresponding to the situation where default occurs the first time that firm value is less than the face value of the loan and the bondholders immediately get the current value of the firm. This lattice is constructed by first filling in the bold cells, which are cells where default occurs (as seen from the firm value lattice in Figure 17.3). The value at a bold cell is equal to the corresponding firm value. Other cells are calculated by the standard backward method, but the bold cell values are not updated. For this example, we find that the value of the bond is \$671,520, which is greater than when default can occur only at maturity. Early default can be better for bondholders because they get most of their money back and get it early.

We can consider credit insurance in this context. Suppose that the bondholders want to eliminate the credit risk. They can do this by purchasing a guarantee. At every default node in Figure 17.5, the insurance company will pay the bondholder the difference between the face value of 800 and the current firm value. We can then determine the value of this “protected” bond. We simply change Figure 17.5 by entering 800 in each default (bold) node. If this is done, the value of the bond will leap up to \$696,680.

Of course, the insurance company will not do this for nothing. Usually payment is made with a **credit default swap**. Such swaps are discussed in more detail in Section 17.11, but in essence the insurance is paid by agreeing to a fixed premium payment at the end of each period that did not result in default. Once default occurs or maturity without default is reached, the contract terminates.

In our example, there is a fixed payment at the end of each year. To determine the proper level of payment, we insert the face value in each default node as before. Next we insert a payment (a cost) of amount, say,  $A$ , into each non-bold cell except the initial cell. We then value this lattice using the standard method of risk-neutral discounting. Finally, we adjust the value  $A$  until the value of the protected bond is equal to the value without protection. In our case, the proper value of  $A$  is \$13.13. The corresponding lattice, accounting for the insurance premiums, is shown in Figure 17.6.

Early Default					
671.52	658.19	684.66	723.87	760.98	800.00
<b>758.27</b>		699.90	715.27	760.98	800.00
	<b>574.97</b>		<b>758.27</b>	742.04	800.00
		<b>435.98</b>		<b>574.97</b>	<b>758.27</b>
			<b>435.98</b>	<b>330.59</b>	<b>435.98</b>
					<b>250.67</b>

**FIGURE 17.5 Lattice for early default.** The bold cells are points of default. Risk-neutral probabilities are used to fill in the other cells by a backward process.

CDS in Place					
671.52	620.08	640.99	686.37	735.36	786.87
	<b>800.00</b>	692.70	689.08	735.36	786.87
		<b>800.00</b>	<b>800.00</b>	741.32	786.87
			<b>800.00</b>	<b>800.00</b>	<b>800.00</b>
				<b>800.00</b>	<b>800.00</b>
					<b>800.00</b>

FIGURE 17.6 Credit default swap applied to protect the bond.

## Coupons\*

The representation of a stock paying fixed dividends or a bond paying fixed coupons typically introduces path dependencies. (See Exercise 5 in Chapter 14.) It is possible, however, to evaluate such processes using a standard binary framework. The trick is that forward construction typically is based on multiplication, whereas in backward evaluation values are summed. Also, to evaluate a coupon bond, we keep track of firm equity rather than debt level. Then, when we have valued the initial equity at time 0, we can subtract that from the firm value to obtain the initial value of debt.

**Example 17.6** Figure 17.7 shows an example using the same HiTech firm and its bond, except that now we have included a modest yearly coupon of \$10 (thousand). We first generate the firm value lattice going forward, as in Figure 17.3. We convert the last column to equity by the formula  $E = \max(0, V - F - C)$ , where  $V$  is the firm value,  $F$  is the face value of the bond, and  $C$  is the coupon. For example, the top node in the last column is  $3989 - 800 - 10 = 3179$  (see figure 17.3). Next we evaluate earlier values of equity by going backward. The value of  $E$  at any prior node is  $E_{\text{new}} = \max[0, (qE^u + (1 - q)E^d)/R] - C$ , where, of course,  $E^u$  and  $E^d$  are the upper and lower next-period values of equity, respectively.

When we get back to the starting point, the value is the value of equity there. The value of the debt (the bond) is the difference between total value of \$1,000 (thousands) and the equity. Hence, the bond value is  $\$1,000 - \$389.3 = \$610.7$  (as compared to

Equity Lattice with Coupon					
389.30	644.91	1018.34	1541.27	2244.44	3179.28
	173.19	327.23	577.03	958.74	1483.70
		45.10	110.82	242.99	508.80
			0.00	0.00	0.00
				0.00	0.00
					0.00

$$\text{Debt value} = 1000 - 389.30 = 610.70$$

**FIGURE 17.7 Coupon bond value.** This method can be used to value a variety of bonds. In this particular case we see that having a coupon increases value but that the default early arrangement is still more valuable.

\$564.34 without the coupon). This technique can be used to explore other issues as well. (See the references.)

## 17.3 Rating Methods

A number of firms publish credit ratings of financial instruments, including government entities and some bond offerings. The ratings are defined by categories, with triple A (expressed as AAA by Standard & Poors and Aaa by Moody's) being the best rating, denoting the highest level of safety. A rating of C or D denotes near or actual default. Figure 17.8 shows the rating categories used by several rating firms.

The classic Merton model forms the basis of some important commercial rating agency methodologies. One such method is that of the KMV corporation (owned by the Moody's corporation), which follows the Merton method fairly closely. The details of the procedure are proprietary, but roughly it models the value of the firm as following geometric Brownian motion. From this, the probability of default can be computed using the lognormal nature of terminal value. This probability is a simple function of the so-called **distance to default**, defined as the number of standard deviations that the logarithm of the firm's value is from the logarithm of the face value of the debt. However, rather than use this value alone, KMV employs historical default data to produce revised probabilities of default, which are then translated into rating categories. An alternative method is that of CreditMetrics. This method also uses conventional rating categories. However, additionally, migrations between categories over a 1-year period are expressed in terms of transition probabilities. One way to obtain these probabilities is based on the Merton model. Alternatively, and

Credit risk	Moody's	Standard & Poors	Fitch Ratings
<b>INVESTMENT GRADE</b>			
Highest quality	Aaa	AAA	AAA
High quality	Aa	AA	AA
Upper medium grade	A	A	A
Medium grade	Baa	BBB	BBB
<b>NOT INVESTMENT GRADE</b>			
Lower medium	Ba	BB	BB
Low grade	B	B	B
Poor quality	Caa	CCC	CCC
Most speculative	Ca	CC	CC
May default	C	C	C
In default	D	D	D

**FIGURE 17.8 Rating categories of major rating companies.** There are slight variations in category names.

From/to	AAA	AA	A	BBB	BB	B	CCC/C	D	NR
AAA	68.89	31.11	0	0	0	0	0	0	0
AA	0	84.09	12.88	0	0	0	0	0	3.03
A	0	0.57	94.26	3.06	0	0	0	0	2.1
BBB	0	0	3.07	91.9	1.54	0.14	0	0	3.35
BB	0	0	0	4.74	82.62	5.19	0.23	0	7.22
B	0	0	0.13	0	6.32	82.06	1.94	1.16	8.39
CCC/C	0	0	0	0	0.63	28.93	37.74	22.64	10.06

**FIGURE 17.9 2010 Corporate Transition Rates: U.S.(%)**

Source: 2010 Annual U.S. Default Study and Rating Transitions, Standard & Poors. 30 March 2011.

more commonly, the 1-year-category transition probabilities are determined from historical records of such changes.

The migration probabilities can be arranged in a transition matrix, an example of which is shown in Figure 17.9. The probability of default at any future year (say,  $n$  years from now) can be found by repeated multiplications of the matrix. This can be done by straightforward computation or by simulation, whereby, starting at an initial category, a jump to a succeeding category is determined by the transition probabilities. Then, after  $n$  steps, the final category is attained. (See Exercise 4.) The average of those that fall in the default category is the estimate of the default probability from the initial point.

This procedure can be generalized in several ways. For example, the rating transitions can be applied separately to the components of a portfolio of securities, although account should be made of the correlation between different securities. As another example, the migration probabilities can be used to determine a credit VaR, which is like a VaR but used to measure credit risk.

## 17.4 Intensity (Reduced-Form) Model

An alternative method for characterizing credit risk does not attempt to model the underlying mechanism of default, as structural methods do, but, rather, characterizes the credit spread that exists in the market for a particular entity—as revealed by current prices of a variety of instruments issued by that entity or similar ones. This method is similar to the use of a short rate process to characterize interest rate risk, where, there too, the changes are not related to an observable variable. This newer approach to credit risk is often referred to as the **intensity method**, because it models the intensity of the predilection to default. The method is alternatively termed the **reduced-form** approach because it does not use the structure of the entity to determine risk.

### Poisson Processes

The intensity model of default assumes that defaults occur randomly but essentially without warning. This type of process, termed a **Poisson process**, is often used in other

contexts to describe the time of occurrence of physical events, such as exhaustion of a light bulb, radioactive decay, and arrivals of visitors to an Internet site. The basic process is defined by a single parameter  $\lambda > 0$ , termed the **intensity**, that quantifies the likelihood of occurrence of an event in a short period. Specifically, the probability of an event's occurrence in a small interval of width  $\Delta t$  is

$$P[\text{event} \in [t, t + \Delta t)] = \lambda \Delta t + o(\Delta t).$$

The probability that two or more events will occur in the interval is  $o(\Delta t)$ ; that is, it goes to zero faster than  $\Delta t$ .

In the context of credit analysis, such a process is meant to describe the random occurrences of specified credit events, such as defaults. Although the basic process continues indefinitely, most credit default models assume that the process is stopped at the time  $\tau$  of the first credit event or at the maturity time  $T$  of the financial entity, whichever comes first.<sup>3</sup>

When this model is used to price a credit-sensitive instrument, the probabilities associated with the model are assumed to be risk neutral. Real probabilities may also be used to measure various statistical properties, such as actual probabilities of default. Most of the time we will be interested in risk-neutral probabilities and use the notion  $q$  rather than  $p$ .

It is important to know the probability that no event occurs in a period of length  $t$ —the **survival probability**—which is easily found. Partitioning the interval  $[0, t]$  into  $m$  segments of length  $\Delta t = t/m$ , the probability of no event in such a segment is 1 minus the probability of an event. Thus the survival probability over the short interval  $t/m$  is approximately  $1 - \lambda t/m$ . Accordingly, the probability of no event over the full interval  $[0, t]$  is approximately  $(1 - \lambda t/m)^m$ . In the limit as  $m \rightarrow \infty$ , the probability  $q(t)$  of no event in time range  $[0, t]$  becomes<sup>4</sup>

$$q(t) = e^{-\lambda t}. \quad (17.8)$$

It follows that the risk-neutral probability that an event does occur at some time  $\tau$  within  $[0, t]$  is

$$\hat{P}[\tau \leq t] = 1 - e^{-\lambda t}.$$

To see how the reduced-form model is used for pricing, consider a zero-coupon bond with face value  $F$  and maturity  $T$  that is not subject to default. Assume that the interest rate  $r$  is constant. Then we know that the value of the bond is

$$B = e^{-rT} F.$$

Now consider a bond that is identical but is subject to default and will in fact default if the associated Poisson process has an event within  $[0, T]$ . If it defaults, the bondholders receive nothing. The value of this bond is simply the discounted risk-neutral expected value

$$V = e^{-rT} \hat{E}[\text{payoff}] = e^{-rT} q(T)F = e^{-rT} e^{-\lambda T} F = e^{-(r+\lambda)T} F,$$

<sup>3</sup> However, if one considers a basket of instruments (say, several bonds), then we may consider several events from the intensity process.

<sup>4</sup> In general,  $\lim_{m \rightarrow \infty} (1 - x/m)^m = e^{-x}$ .

where  $q(T)$  is the risk-neutral survival probability. This shows that  $\lambda$  acts like an increment to the interest rate; hence, it is termed the **short-term credit spread**.

## Inhomogeneous Process

Now suppose that the intensity  $\lambda$  varies with time according to a function  $\lambda(t)$ . The resulting Poisson process is then said to be **inhomogeneous**, as compared to a **homogenous** process, in which  $\lambda$  is constant. To find the survival probability of the inhomogeneous model in the interval  $[0, t]$  we again partition  $[0, t]$  into  $m$  segments of length  $\Delta s = t/m$  and let  $s_i = i\Delta s$ . The probability of zero events over the  $i$ th interval is approximately  $e^{-\lambda(s_i)\Delta s}$ . The survival probability  $q(t)$  is then

$$\begin{aligned} q(t) &= \lim_{\Delta s \rightarrow 0} (e^{-\lambda(s_1)\Delta s})(e^{-\lambda(s_2)\Delta s}) \dots (e^{-\lambda(s_m)\Delta s}) \\ &= \lim_{\Delta s \rightarrow 0} e^{\sum_1^m \lambda(s_i)\Delta s} = e^{-\int_0^t \lambda(s)ds}. \end{aligned} \quad (17.9)$$

Suppose that the short rate of interest is known to follow a curve  $r(t)$ . The price of a defaultable zero-coupon bond with face value  $F$  that pays nothing in the event of default is again the discounted risk-neutral expected value of the payoff. In this case

$$V = \hat{E}\left[e^{-\int_0^T r(s)ds} F\right] = e^{-\int_0^T r(s)ds} q(T)F = e^{-\int_0^T (r(s)+\lambda(s))ds} F.$$

The intensity curve defines a credit spread curve directly.

## 17.5 Stochastic Intensity Model<sup>\*</sup>

In reality, an intensity curve  $\lambda(t)$  is likely to vary randomly with time in a manner similar to a short rate process for interest. There are many ways this could be represented, but for analytic and computational tractability it is important that the randomness be generated in two steps. First, from a stochastic risk-neutral intensity process (such as the CIR process to be discussed in Section 17.7) a specific path of intensities  $\{\lambda(s) : 0 \leq s \leq T\}$  is generated. Next, given that path, the default time  $\tau$  is generated according to the inhomogeneous Poisson process with the given path intensities. The result is termed a **doubly stochastic** or **Cox** process. The survival probability of this process is found by going backward in the two-step procedure. First, the survival probability corresponding to the specific path of the inhomogeneous process is found as  $e^{-\int_0^T \lambda(s)ds}$ . Next, based on the random process for  $\lambda(t)$ , the risk-neutral expected value of this is evaluated. Thus

$$\hat{P}[(\text{event} > T) | \lambda(s) : 0 \leq s \leq T] = [e^{-\int_0^T \lambda(s)ds}],$$

as derived in Section 17.4. We have the simple formula

$$q(T) = \hat{P}(\text{event} > T) = \hat{E}[e^{-\int_0^T \lambda(s)ds}]. \quad (17.10)$$

The time fluctuations of  $\lambda$  can be represented by many of the methods that have been used in earlier chapters, including an Ito process, a discrete-time process that

can be used in simulation, or a lattice model. However, the  $\lambda$ 's must be nonnegative. When we put all of this together, the processes and methods for evaluating credit derivatives are similar to the methods used to evaluate interest rate derivatives using an underlying short rate process. We simply take the discounted risk-neutral expected value with respect to all possible (risk-neutral) realizations of the processes.

Specifically, for a zero-coupon bond with face value  $F$ , the value of the bond at time 0 is

$$V = \hat{E} \left[ e^{-\int_0^T [r(u) + \lambda(u)] du} F \right]. \quad (17.11)$$

This can be interpreted as the risk-neutral average over paths of  $\lambda$ .

## 17.6 Intermediate Receipts

So far we have considered only the case where, if default occurs, no payment is made to bondholders. In reality, bondholders may receive some portion of the face value. If the recovery payment is of known fixed amount and is received at maturity, it is easy to evaluate the bond. For example, if the recovery payment is  $c$ , we recognize that the risk-neutral expected value of the total payment is

$$q(T)F + [1 - q(T)]c.$$

Actually, it is likely that some recovery payment will be received at or shortly after default. To assess such situations (and to value various other derivatives of credit-sensitive instruments), it is necessary to value cash flows that occur at the time of default. This value is governed by the probability density of the (random) default time  $\tau$ .

In the inhomogeneous case, the density at  $t$  is the rate of change of the total probability of default. That is, the density is

$$\begin{aligned} h(t) &= \frac{d}{dt} \text{Prob}[\tau \leq t] = \frac{d}{dt} (1 - q(t)) = -\frac{d}{dt} q(t) \\ &= \lambda(t) e^{-\int_0^t \lambda(s) ds}. \end{aligned} \quad (17.12)$$

In the special case of a homogeneous process, this becomes

$$h(t) = \lambda e^{-\lambda t}. \quad (17.13)$$

In the general doubly stochastic model, the value of receiving \$1 at the time of default is

$$V = \int_0^T \hat{E} \left[ e^{-\int_0^u [r(s) + \lambda(s)] ds} \lambda(u) \right] du.$$

**Example 17.7 (One dollar)** Suppose that both  $r$  and  $\lambda$  are constant. The default time density is  $h(t) = \lambda e^{-\lambda t}$  and the value of \$1 to be paid at the time of default is

$$D = \int_0^T \lambda e^{-(r+\lambda)u} du = \frac{\lambda}{r+\lambda} [1 - e^{-(r+\lambda)T}]. \quad (17.14)$$

For instance, consider a 10-year zero-coupon bond when the interest rate is 5% and the credit spread is  $\lambda = 1\%$ . The value of a dollar to be paid on default is then  $D = \frac{.01}{.06}[1 - e^{-6}] = \frac{1}{6}[1 - .549] = .075$ , that is,  $7\frac{1}{2}$  cents. By comparison, a dollar to be received at maturity in the event of default is worth  $(1 - e^{-1}) \exp(-.5) = .0577$ , that is, about 5.8 cents.

## 17.7 Analytically Tractable Cox Processes

Models used to generate  $\lambda$  in reduced-form methods for credit risk are typically direct analogies to the models used to generate short rates in an interest rate model. Often, an affine model (as discussed in Sections 16.10 and 16.11) is used.

A popular process is the Cox, Ingersoll, and Ross (CIR) model

$$d\lambda = c(\mu - \lambda) dt + \sigma \sqrt{\lambda} dz.$$

For  $2c\mu > \sigma^2$  and  $\lambda(0) > 0$ , the process remains strictly positive. Another important property is that it has a degree of mean reversion. If  $\lambda$  should fall below  $\mu$ , there is a tendency for  $\lambda$  to increase; conversely, if  $\lambda$  climbs above  $\mu$ , there is a tendency for  $\lambda$  to decrease.

Furthermore, the survival probability is

$$q(t) = \exp[a(t) - b(t)\lambda(t)], \quad (17.15)$$

where

$$\begin{aligned} b(t) &= \frac{2(e^{\gamma t} - 1)}{(\gamma - c)(e^{\gamma t} - 1) + 2\gamma} \\ a(t) &= \frac{2c\mu}{\sigma^2} \log\left(\frac{2\gamma e^{(\gamma - c)t/2}}{(\gamma - c)(e^{\gamma t} - 1) + 2\gamma}\right) \end{aligned}$$

and  $\gamma = \sqrt{c^2 + 2\sigma^2}$ . This explicit formula is analogous to the zero-coupon price in the CIR setting. The formula for the survival probability provides a means for explicit evaluation of several credit derivatives.

In basic form, these models are all single-factor models, since they are driven by a single stochastic source. If any one of these is combined with a single-factor model for short rates of interest, the overall model becomes a two-factor model. Other factors can be incorporated as well (such as a process for stock market fluctuations), resulting in even higher-order models.

### Model Fitting

In the early period of formal risk management, the number of credit-sensitive instruments on a single entity was relatively low. A major firm typically had only a few outstanding bonds. Likewise, there were few related credit protection instruments. In such a world it makes sense to model the issuing firm, perhaps as a stochastic process, and to use that model to estimate the likelihood of default. However,

in a world in which credit risk is more acute—where major firms typically have several outstanding bond issues and there are perhaps dozens of secondary instruments in a broad and active market—it perhaps makes sense to characterize the market rather than the firm. The reduced-form approach is of this form. Its general outline is a focus on the credit spread, and its parameters are adjusted by fitting it to the market, using the many existing derivatives. It is analogous to the analysis of interest rate derivatives by postulating a risk-neutral framework and then fitting that framework to match closely the term structure. In credit risk, one fits a spread curve.

**Example 17.8 (Simple fit)** Suppose the price of a 5-year zero-coupon bond is 72.25. The risk-free rate is  $r = 6\%$ . A new bond offered by the same entity is a 3-year zero-coupon bond. What is the expected price of this bond? As a simple procedure, we assume that the intensity for this entity is constant. We estimate  $\lambda$  by writing the price of the existing 5-year bond as  $72.25 = e^{-5(0.06+\lambda)} 100$ . Solving, we find  $\lambda = .005$ . Using this estimate, we find the price of the 3-year bond as  $e^{-3(0.06+.005)} = 82.28$ .

## 17.8 Simulation

Complex financial instruments with special procedures for restitution, having periodic cash flows, and represented by multiple-factor stochastic intensity models are a challenge for an analyst. It is generally not possible to obtain closed-form expressions for appropriate prices in these cases. Simulation then becomes a natural alternative; accordingly we outline a few basic methods that use Monte Carlo simulation.

The general Monte Carlo method of simulation was described in Section 15.8. Basically, to evaluate the expectation  $E[f(X)]$  when  $X$  is a random variable, one generates a series of independent samples of  $X$ , say,  $x_i$ ,  $i = 1, 2, \dots, N$ , and approximates the expected value as

$$E[f(X)] \approx \frac{1}{N} \sum_{i=1}^N f(x_i).$$

The standard deviation of the error goes down according to a factor times  $1/\sqrt{N}$ ; hence, in practice the method typically requires a very large  $N$ . This is especially true if  $X$  is multidimensional and high accuracy is required in the tail of the distribution. The methods described next use Monte Carlo, but in somewhat different ways.

### Direct Simulation

Roughly, in the direct approach, the process is carried out in discrete time steps of length  $\Delta t$ . At each step, the risk-neutral  $\lambda$  (along with  $r$  and any other factors) is updated according to the underlying process. The probability of a default during the step is  $1 - e^{-\lambda \Delta t}$ . Accordingly, we generate a random variable with uniform density on  $[0, 1]$ . If the result (say,  $u$ ) satisfies  $u \leq 1 - e^{-\lambda \Delta t}$ , we declare that the process

defaults that period; otherwise not, and the process continues to the next step.<sup>5</sup> If there is a default, we calculate the payment that would be due and discount it back to the initial time, according to the current risk-neutral discount factor. A single run of the process is terminated at default or at maturity, whichever comes first.

The method is flexible and straightforward. Unfortunately, it is also terribly slow. The main cause of the sluggishness is the determination of the stopping times, requiring that a uniform random variable be generated at each step. Furthermore, since the probability of default is small yet has a significant impact on value, it requires many, many trials to obtain reasonable accuracy. See Exercise 11.

## A Better Way

In essence, the direct method generates thousands (or millions) of paths and corresponding default events, each time recording the appropriate payoffs, and then averages the results. The end result is the overall value: the risk-neutral expected value of the discounted payoffs. We can speed up the process dramatically by averaging as we go.

For example, as the process moves to the next period of length  $\Delta t$ , the direct method makes a decision as to whether to declare a default or not, according to the risk-neutral probability of default; then records the resulting payoffs (discounted for interest); and finally moves ahead unless there was a default.

In an alternative approach, when the process moves to a new period, the payoffs are evaluated each way—default and survival—and these are weighted by their risk-neutral probabilities. The resulting payoffs are recorded by interest rate discounting and by discounting according to survival to that point. The process moves on in all cases.

Specifically, the overall simulation process consists of the following steps:

1. Select a time step  $\Delta t$  and a model that can generate the short rate and intensity values at each time step.
2. Move through the process, step by step, and at each step update  $r, \lambda$ , and any other random factors to their current values.
3. Let  $q = e^{-\lambda \Delta t}$  be the survival probability for a given step. Calculate all current cash flows  $c_d$  and  $c_s$  for default and survival, respectively, in this period, given that there has been survival to this point. Let  $C = (1 - q)c_d + qc_s$ . Multiply  $C$  by the current survival probability (from  $t = 0$ ) and discount back to  $t = 0$  to account for interest.
4. Sum all discounted cash flows to produce the value associated with this trial.
5. Compute the average over many trials to obtain an estimate of the total value.

Note that a default is never generated in this evaluation process. Instead, both possibilities are accounted for, in the sense of their average influence. The method is very

---

<sup>5</sup> See Exercise 9 for a somewhat different way.

powerful, but in basic form it cannot treat efficiently complexities such as path dependency, whereas the direct method can. However, neither method can treat situations that have decision opportunities.

## 17.9 Lattice Methods

Lattices can be quite effective for evaluation of credit derivatives. They are especially useful for situations where there is opportunity for decisions, such as for credit default swap options or callable bonds. Lattices also have the feature that, like analytical solutions, once constructed, they almost instantly produce revised results in response to changes in parameter values. Their disadvantage is, of course, that they may be very large when representing models with several factors or for a path-dependent situation that requires a tree rather than a lattice. However, the study of lattice methods is quite useful for understanding clearly the underlying pricing process.

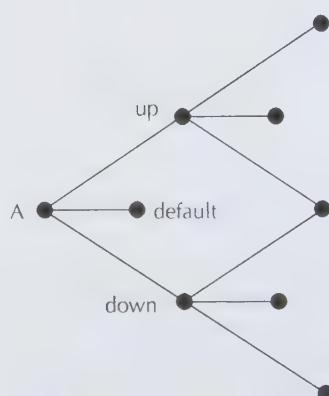
A standard reduced-form model for valuing credit derivatives incorporates two factors: the short rate process for interest, and the stochastic model for intensities. Accordingly, a suitable lattice will generally have dimension of at least 2. A good approach is to begin by building two separate lattices—one for  $\lambda$  and one for  $r$ —and then integrate them appropriately to match observed market prices, interest rates, volatilities, and covariances.

We shall outline a method that uses binomial lattices. We focus primarily on the intensity lattice and, for simplicity, assume that the short rate of interest is constant. Later we will explain how to expand the model to include a short rate lattice.

First, a specific period length  $\Delta t$  is established. The lattice grows forward period by period. At any point A in the lattice, it is possible to move forward along any one of *three* branches: “up,” “default,” or “down” as shown in Figure 17.10. The default branch from A has a risk-neutral probability determined by the intensity at A; specifically,  $q_{\text{default}} = 1 - \exp(-\lambda \Delta t)$ , where  $\lambda$  is the current (node A) intensity.

It is common in an intensity model to include a degree of mean reversion. In a lattice, this feature can be included by assigning different risk-neutral probabilities to different branches, depending on  $\lambda$ . We propose the formulas for the “up” and

**FIGURE 17.10 Beginning of an intensity lattice.** There are 3 possible successor nodes, but the default node is a stopping point for the process.



“down” probabilities as

$$q_{\text{up}} = \frac{\lambda_0}{\lambda + \lambda_0} - \frac{1}{2}q_{\text{default}} \quad (17.16a)$$

$$q_{\text{down}} = 1 - \frac{\lambda_0}{\lambda + \lambda_0} - \frac{1}{2}q_{\text{default}}, \quad (17.16b)$$

where  $\lambda_0$  is a reference value. Clearly,  $q_{\text{down}} + q_{\text{default}} + q_{\text{up}} = 1$ . Notice that equations (17.16) imply that the up and down probabilities are equal when  $\lambda = \lambda_0$ . If  $\lambda$  should move upward from this, then the probability of an up move decreases, inducing a tendency for  $\lambda$  to move back toward  $\lambda_0$ . The reverse occurs when  $\lambda$  falls below  $\lambda_0$ . Of course, alternatives to the simple structure (17.16) can be postulated, and even this one could be easily varied by endowing it with a variable speed of attraction<sup>6</sup> to  $\lambda_0$ .

The lattice is completed by assigning specific  $\lambda$  values to the lattice nodes. We do this by using a procedure similar to the Ho–Lee method for short rates described in Section 16.7 in Chapter 16, as first proposed by Black, Derman, and Toy. As in that method, the nodes are labeled  $(k, s)$ , where  $k$  and  $s$  are the time and state indices, respectively. The lattice is populated by values of  $\lambda$  indexed as  $\lambda_{k,s}$ . We specify the form

$$\lambda_{k,s} = a_k + bs, \quad (17.17)$$

where  $b$  is a volatility parameter.

**Example 17.9 (Please buy Junk)** Mr. Julio Junk has built a multimillion-dollar import–export business in rare antiques. His success was fueled by obtaining financing through a series of bond offerings. He now offers a \$100 million, 10-year, 6% coupon bond and is trying to convince Major Bank to buy a large piece of it. People at the bank are leary of investing in Junk’s bonds, but they are willing to consider it. First they want to price it. They ask us to build a lattice assuming that there will be no recovery if there is default.

We construct the lattice for a 10-year period, with a fixed interest rate of 5% and a period length of  $\Delta t = 1$  year. The lattice entries are determined from equation (17.17) with  $b = .002$  and a set of  $a_k$ ’s that will be fitted to the observed term structure of credit (obtained in practice through study of existing prices of other bonds of that entity). The yearly short rates for credit are  $Y(t) = .08 - .02 \exp(-t/2)$  and are shown in the top line of Figure 17.11. The probabilities are defined by the reversion to the mean formula (17.16) with  $\lambda_0 = 2\%$ .

The completed lattice is shown in the top part of Figure 17.11. The corresponding  $q$  lattice of risk-neutral “up” probabilities is shown only to illustrate the degree of mean reversion. (The up probability decreases with node level, which tends to exert a downward force.) The main lambda lattice is constructed by selecting a set of  $a_k$ ’s (along the bottom row of the  $\lambda$  lattice), so that the corresponding short rates match those specified. This matching process is identical to that used in Chapter 16 for the

<sup>6</sup> For example, we could set  $q_{\text{up}} = (\lambda_0 / (\lambda + \lambda_0))^\gamma - \frac{1}{2}q_{\text{default}}$  for a suitable constant  $\gamma$ .

	0	1	2	3	4	5	6	7	8	9	10
	Credit Short Rates										
0.060	0.068	0.073	0.076	0.077	0.078	0.079	0.079	0.080	0.080	0.080	
											0.043
											0.042
											0.041
											0.041
											0.039
											0.038
											0.037
<b>Lambda Lattice</b>						0.040	0.038	0.037	0.036	0.035	
						0.039	0.038	0.036	0.035	0.034	0.033
						0.039	0.037	0.036	0.034	0.033	0.031
						0.037	0.037	0.035	0.034	0.032	0.029
						0.034	0.035	0.033	0.032	0.030	0.027
						0.026	0.032	0.033	0.031	0.030	0.025
						0.010	0.024	0.030	0.031	0.029	0.023
<b>q Lattice</b>						0.315	0.325	0.333	0.342	0.350	
						0.318	0.328	0.338	0.347	0.356	0.365
						0.322	0.331	0.341	0.352	0.362	0.372
						0.331	0.335	0.345	0.356	0.367	0.378
						0.354	0.345	0.349	0.359	0.371	0.383
						0.418	0.369	0.359	0.364	0.375	0.387
						0.662	0.438	0.385	0.375	0.380	0.392
											0.299
											0.304
											0.310
											0.323
											0.328
											0.336
											0.350
											0.365
											0.381
											0.398
											0.416
											0.437
											0.459

**FIGURE 17.11** The top lattice is the completed intensity lattice matched to the given term structure for credit. The lower lattice shows the risk-neutral probabilities  $q_{up}$ , illustrating mean reversion.

Ho–Lee short rate lattice of  $r$ 's. The  $a_k$ 's in the lattice of Figure 17.11 are the result of this match.<sup>7</sup>

The main lattice is now available for analysis of securities that are derivatives of the credit risk of this entity. The short rate curve of risk-free rates is flat at 5% out to 10 years.

The value of Junk's bond, accounting for the possibility of default with zero recovery, is found to be \$66.79 million from Figure 17.12. As expected, this is well below the value of \$106.7 million that would be the price if there were absolutely no credit risk.

The node values in the figure consist of three parts that are summed. (1) The risk-neutral probability of default that period,  $1 - \exp(-\lambda \Delta t)$ , is multiplied by the

<sup>7</sup> Throughout we chose to discount at the continuously compounded rate each year. That is, we discounted by  $\exp(-r \Delta t)$  instead of  $1/(1+r \Delta t)$ .

0	1	2	3	4	5	6	7	8	9	10
										101.58
									94.78	101.79
							89.16	95.33	101.99	
						84.50	89.98	95.87	102.19	
					Valuation Basic	80.63	85.55	90.81	96.42	102.40
				77.39		81.86	86.61	91.65	96.98	102.60
			74.73		78.78	83.11	87.68	92.49	97.54	102.81
		72.64		73.24	80.19	84.38	88.77	93.35	98.10	103.02
	71.26		74.25	77.78	81.62	85.66	89.87	94.20	98.66	103.22
71.03		72.98	75.91	79.35	83.08	86.97	90.98	95.07	99.23	103.43
66.79	72.85	74.74	77.60	80.96	84.56	88.30	92.10	95.94	99.80	103.64

**FIGURE 17.12 Basic valuation of Julio Junk's bond.** The short rate of interest is flat at 5%. The credit spread is defined by the intensity process of Figure 17.11.

amount received there at default (which is zero in this case). (2) Next, the risk-neutral survival probability is multiplied by the coupon payment. (3) Finally, the risk-neutral probability of survival is multiplied by the usual up and down valuation of next-period's value, and this is discounted by one period at the current short rate of interest.

Specifically, the values in the final column are found as the survival probability for one period times the face value plus the coupon. Hence, for the top node,  $e^{-.043} 106 = 101.58$ . The value at the top of the next column is

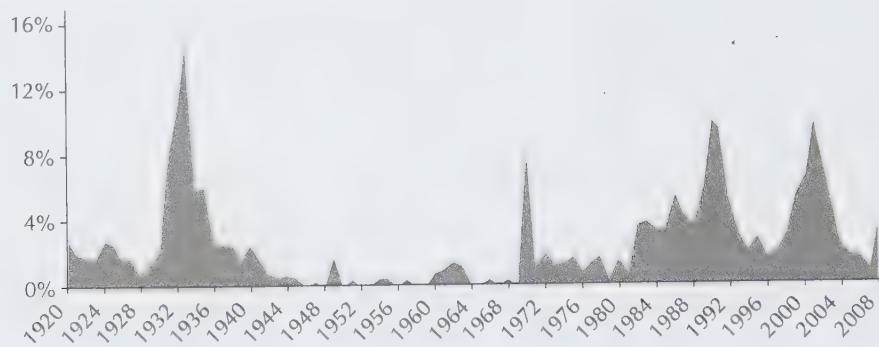
$$e^{-.042} \{6 + e^{-.05} [q_{\text{up}} 101.58 + q_{\text{down}} 101.79]\} = 94.78.$$

In Section 17.11 the model is used to evaluate a credit default swap on Junk's bond.

In general, it is necessary to include a lattice representing the term structure of interest rates. This  $r$  lattice can be constructed first, independent of the intensity lattice, and made to match the existing term structure of interest rates. Next, the intensity lattice is constructed by incorporating the interest rate lattice. The final lattice is a composite of these two, and this master lattice will have a total of five successor nodes from each point. If  $\lambda$  and  $r$  are (risk-neutral) independent, then the probabilities of the nodes in the master lattice are simple products. If  $\lambda$  and  $r$  are correlated, then the proper values can be found by the method described in Section 19.5 in Chapter 19.

## 17.10 Correlated Defaults

A history of actual defaults in the United States shows that defaults do not arise independently over time but, rather, occur in clusters. See Figure 17.13. These clusters tend to include several companies from the same industry. For example, 22 oil companies defaulted between 1982 and 1986, and 747 savings and loan institutions failed



**FIGURE 17.13 Annual speculative-grade default rates 1920–2008.** Source: Moody's Default Risk Service.

in the 1980s and 1990s crises. Several modifications or extensions of the basic credit models have been proposed to account for this phenomenon.

A group of  $n$  credit-sensitive entities may play two distinct investment roles. First, they may constitute a portfolio of assets in which different amounts of each may be taken, as in the standard portfolio problem. The overall credit risk is determined by the weights and the credit correlations among entities. Second, these  $n$  entities may constitute a fixed basket for which a master credit insurance instrument is designed, with payoff determined by given conditions, such as whenever there is a default of one of the constituent entities or when there has been a specified number of defaults (such as three) within the basket. In either case the  $n$  entities may be known to be interdependent; hence, information about one may affect the probability of default by others.

It is worth noting that although the term **correlation** is used to describe interdependence among entities, the term is not necessarily meant in the restricted sense of standard (so-called linear) correlation. Some forms of dependence may not be well characterized by this standard definition. For example, suppose one firm in an industry is a leader and another a follower. If the leader fails, it is likely that the follower will also fail. But the reverse may not hold.

Most methods for analyzing credit risk in the presence of correlation are extensions of methods used for evaluating a single entity. Largely, they fall into the following four categories.

1. **Credit rating methods.** It is natural to extend single entity migration of credit ratings to joint transitions in order to account for correlation. An expanded set of transition probabilities governing the joint movement of several correlated entities is required, and clearly the dimension of such a transition matrix can be very large compared to one that did not consider correlation. Estimation of the entries of the transition matrix can be difficult.
2. **Structural methods.** The underlying model of firm value used in the structural approach can be generalized to include correlation among the values of several

firms. Then the probability of default of one entity triggered by another company's value falling below a boundary can be calculated (most likely by simulation). Also, the introduction of processes that are jumpy in character can sometimes mimic realistic behavior.

3. **Reduced-form methods.** The stochastic differential equation generating the intensity  $\lambda$  can be generalized to a set of processes governing the  $n$  intensities. Or a single process may be sufficient by considering a sequence of events. Copulas, defined next, represent a popular reduced-form approach.
4. **Copulas.** For the purpose of simulation, it is useful to generate a set of  $n$  suitably correlated default times  $\tau_1, \tau_2, \dots, \tau_n$ . Although they are correlated, we wish that, individually, each  $\tau_i$  should have the probability distribution implied by its own survival probability function  $p_i(t_i)$ . (Here we use  $p$  to indicate actual probability, but the same process applies to risk-neutral probabilities.) We introduce the correlation by what is known as a **copula**. There is a general copula theory, but for this application the most popular, by far, is the Gaussian copula, and that is what we shall discuss. There are three steps.
  - (a) Generate a set of  $n$  standard normal (Gaussian) random variables  $X_1, X_2, \dots, X_n$  with a given correlation matrix. This is relatively easy.
  - (b) Define the associated random variables  $U_i = F_N(X_i)$ , where  $F_N$  is the distribution function of a standard normal. Note that each  $U_i$  is uniformly distributed on  $[0, 1]$  because

$$P[U_i \leq u_i] = P[F_N(X_i) \leq u_i] = P[X_i \leq F_N^{-1}(u_i)] = u_i.$$

Although they are uniform, the  $U_i$ 's are generally not independent because the underlying  $X_i$ 's are correlated.

- (c) For each  $i$  let  $\tau_i$  satisfy  $p_i(\tau_i) = U_i$ , where  $p_i$  is the  $i$ th survival function. We have

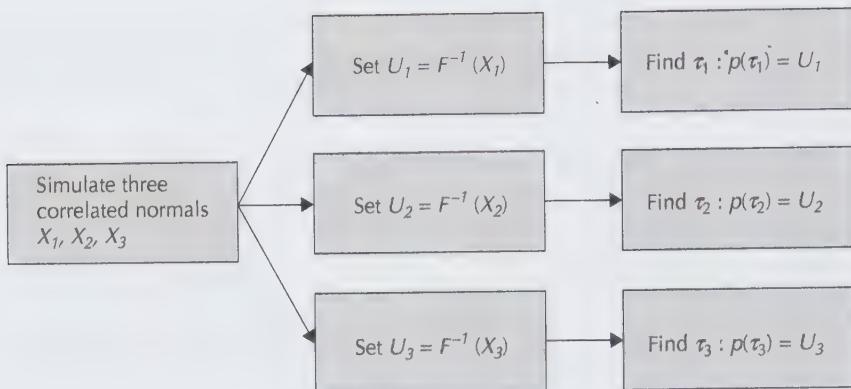
$$P[\tau_i \geq t_i] = P[p_i^{-1}(U_i) \geq t_i] = P[U_i \geq p_i(t_i)] = p_i(t_i),$$

the last step holding because  $U_i$  is uniform. Thus, the result is a set of random default times that are consistent with the known individual survival functions and that exhibit correlation transformed from the joint normal distribution.

The process is illustrated in Figure 17.14 for the case of three default times.

## 17.11 Credit Derivatives

Cataloging the various forms of credit derivatives that have been devised in just the past few years would be an enormous task, and the resulting large document would have to be updated frequently as new products emerge. Therefore, we offer here only a sampling of the prominent forms of such derivatives. Also a few illustrative analyses are presented using simple or slightly complex methods, to provide a flavor of what can be done.



**FIGURE 17.14 Simulation using a Gaussian copula for three variables.** First, three jointly standard normal random variables are generated with a specific covariance structure. From these, three uniform variables on  $[0, 1]$  are defined by applying the inverses of the standard normal distributions. Finally, the stopping times are found through the inverse of the individual survival time distributions.

## Bonds and Loans

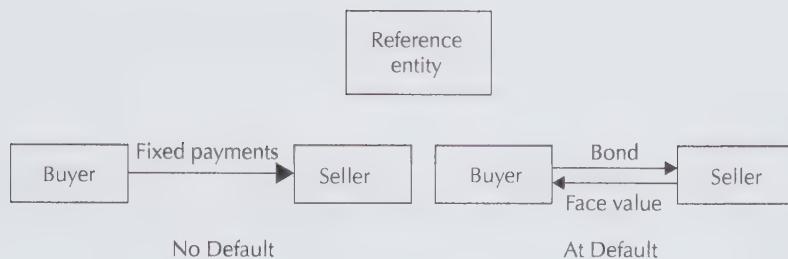
Bonds have been discussed throughout the chapter, since they represent the basic form of a credit-sensitive instrument. In a standard bond, a fixed cash flow is promised, consisting of fixed periodic payments (coupons) together with a face-value payment at maturity. If the bond is held to maturity, there is no risk to the stream of cash flows other than credit risk. The value of the bond varies at other times, depending on the interest rate.

Sometimes there are special features, such as whether a bond is callable or whether it carries with it the option to convert it to stock.

Loans are similar, except that in some cases the required payments are tied to some other variable—most often a specific interest rate—as when a construction loan requires payments tied to (floating) LIBOR.

## Credit Default Swaps (CDS's)

Credit default swaps are the most common derivatives used for protection against default. They are structured much like an insurance policy: Party A, who seeks protection, agrees to pay party B, the insurer, an annuity in return for protection against default on a particular bond (the reference obligation) of entity C. To carry this out, a notional amount is established, say, \$100 million. Party A agrees to make periodic payments to B in equal amounts determined by the **CDS spread**. If this spread is 300 basis points, then A's payments are 3% of the notional value each year. Payments are made quarterly, and they are made in arrears each period, meaning that they are paid at the end of each period of protection. Payments stop at default. If the reference entity C defaults on the reference bond, then, in physical settlement, A has the right to



**FIGURE 17.15 CDS organization.** The buyer of protection makes periodic payments to the seller of protection. On default, in physical settlement the buyer sends the seller the bond of the reference entity and the seller sends the face value to the buyer.

deliver to B the bond of C with face value equal to the notional amount and receive the face value from B. For example, A may own a \$100 million face-value bond issued by C that is to be protected. In case of default, A can exchange this bond for its face value. Alternatively, if A does not own the reference bond, then settlement is made in cash. A receives the difference between the face value and the current market value of the bond. (See Figure 17.15.) If default occurs partway through one of the payment periods, then A must pay the corresponding fraction of that period's payment to B.

In some arrangements a **digital swap** is defined where, on default, B will pay a predetermined amount to A.

One way to structure CDS's is to set the level of fixed annuity payments to satisfy the equation

$$\text{Protected value} - \text{Annuity value} = \text{Value without protection}. \quad (17.18)$$

The protected value is the value assuming that in the event of default full payment of face value will be delivered to the buyer of protection. This value will generally not be equal to the value of an identical bond that is not subject to default.

Equation (17.18) can be written in a slightly different way as

$$\text{PV(premium payments)} = \text{PV(protection payments)}.$$

**Example 17.10 (Swap Junk now)** The people at Major Bank are still quite hesitant to buy a significant portion of Junk's bond. (See Example 17.9.) However, they might consider a purchase if they could also buy a corresponding credit default swap. Again we are asked to help.

We find the proper swap fee (the annuity payments) in two steps. First, the value under full recovery at time of default is found, simply by inserting the face value of the bond (but not the coupon) at every default node. This will change the value of the bond from \$66.79 million to \$86.21 million. This is the value with complete protection from default. To find the proper swap fee, we adjust the coupon downward until the value of this modified bond is equal to the value without protection. The appropriate coupon for Junk's bond is 2.61%, which is equivalent to getting the original 6%

	0	1	2	3	4	5	6	7	8	9	10
<b>Protected Valuation</b>											102.50
											96.29 102.51
											91.02 96.47 102.51
							86.52	91.33	96.65	102.52	
							82.65	86.92	91.64	96.83	102.52
						79.32	83.13	87.33	91.95	97.01	102.53
					76.45	79.84	83.61	87.74	92.26	97.19	102.53
			73.99	77.00	80.37	84.09	88.16	92.58	97.37	102.54	
		71.95	74.57	77.57	80.92	84.59	88.58	92.90	97.56	102.54	
	70.37	72.54	75.16	78.15	81.47	85.09	89.01	93.23	97.74	102.55	
	66.80	70.98	73.15	75.77	78.75	82.03	85.61	89.45	93.56	97.93	102.55

**FIGURE 17.16 Valuation of Julio Junk's bond.** This is with full protection and a swap fee of 339 basis points.

coupon and paying 3.39% as a swap payment. Hence, the swap requires a payment of 339 basis points.

The lattice that corresponds to full protection and a fee of 339 basis points is shown in Figure 17.16. A typical node is evaluated with the formula  $V = e^{-\lambda} \{C + e^{-r}[q_{\text{up}}V_{\text{up}} + q_{\text{down}}V_{\text{down}}]\} + (1 - e^{-\lambda})100$ , where  $C = 6$  for full protection and  $C = 2.61$  to obtain the lattice shown. Notice that the variance of the numbers in any column is lower than in the original lattice. Of course, the main risk reduction is the protection against default.

## Forwards and Options on CDS's

The CDS market is huge and fairly liquid. It is possible, therefore, to define forwards and options on CDS, much like options or forwards on stocks or commodities.

## Total Return Swaps (TRS's)

Also referred to as TRORS's, for total rate-of-return swaps, these instruments make it possible to hedge market risk as well as credit risk. Suppose that party A owns a bond and seeks protection from party B. A will transfer all coupon payments to B as well as all capital gains (due to interest rate fluctuations). In return, B will pay LIBOR plus a fixed spread on a continuing basis. B will also pay A any capital losses. In the event of default, B pays A the original cost of the bond.

The total return swap can be regarded from the opposite viewpoint as well. Party B may wish to own the income stream from the bond but would like to finance its purchase. Hence, B borrows the cost of the bond from A and pays LIBOR plus a spread for that loan. It then receives the income and capital gains from A, as if B owned the bond.

## Collateralized Debt Obligations (CDO's)

Here, a number of income-producing assets are bundled and fractions of the bundle are sold to investors. The assets may include almost anything that produces income, including bonds, mortgages, auto loans, CDS's, and combinations of these. Often, the ownership of the bundle (or pool) is structured into tranches, similar to the structure for CMO's discussed in Section 16.9, that determine how the income from the bundle is to be distributed to the investors. For cash CDO's, the underlying assets are cash instruments (generally bonds). Each tranche is promised a certain (nominal) percent return on its investment. For example, the first tranche may be promised 6%, the second 8%, and so forth, based on the assumption of no defaults. The income from the bundle is allocated to the first tranche until its 6% target is met; income then goes to the second tranche until its 8% is met, and so forth through the tranches. Due to possible defaults, not every tranche will receive its promised amount. Clearly, higher tranches carry a greater risk than do lower ones.

**Example 17.11 (Two are not always equal)** Consider a CDO pool consisting of just two bonds, each of face value \$100 million. The bonds are identical in structure, and their risk-neutral defaults are independent, with probabilities each equal to  $q = 10\%$ . Using a simplified analysis we show how, by defining two tranches, the risk characteristics can be changed so that one tranche is less risky and the other is more risky than either would be if treated separately. We label the two tranches A and B. The combined face value is \$200 million. Of this, \$120 million is allocated to A and \$80 million is allocated to B. However, A is considered the senior tranche and is paid first, up to A's promised amount. What remains is paid to B. In the event of a default, the recovery rate is 30%. The risk-free interest rate is 5%. Maturity is in 1 year.

The entire situation and analysis is summarized in the Table 17.1. There are three scenarios corresponding to either zero, one, or two defaults. The risk-neutral

**TABLE 17.1**  
**A CDO WITH TWO TRANCHES**

Defaults	Prob	Total payout	Amount to A	Amount to B
0	.81	200	120	80
1	.18	130	120	10
2	.01	60	60	0
		One	A	B
Mean payout		93.000	119.400	66.600
price		88.464	113.577	63.352
Promised yield		.123	.055	.233
Default prob		.010	.010	.190
Sigma*		21.00	5.970	27.685

A is promised \$120,000, while B is promised \$80,000. However, A has priority over B. As a result, tranche A is more expensive and has lower yield but less variance than tranche B.

\*Note: The variance is computed using the risk-neutral probability rather than the (unspecified) actual probability.

probabilities of these three possibilities are shown in the column on the left (using the independence assumption). The total income available under the three possible scenarios is shown next and its disposition to A and B in the following two columns in the upper part of the table. The analysis is shown in the lower part of the table. The column labeled “one” corresponds to the properties of a single bond.

The price of one bond (if sold separately) is its discounted risk-neutral expected value  $[.9 \times 100 + .1 \times 30] \exp(-.05) = 88.46$ . Likewise, the price of tranche A is  $[.81 \times 120 + .18 \times 120 + .01 \times 60] \exp(-.05) = 113.58$ . The promised yield is the logarithm of the ratio of promised payout to price. For tranche A this is  $\ln(120/113.58) = .055$ .

The prices of the single bond and the two tranches are shown in the table. Notice that the prices for A and B sum up to twice the price of a single bond, as they should. The table shows that tranche A is safer but has a lower yield than both the bond itself and tranche B. The default probability of each case is the risk-neutral probability that the full promised payout will not be received.

**Example 17.12 (Correlated assets)\*** Typically the default propensities of the assets in a CDO bundle may be correlated. We can model this for the previous example by considering default variables  $D_1$  and  $D_2$  for bonds 1 and 2, defined to be either 1 for default or 0 for no default. The risk-neutral probability of a 1 is  $q = .10$  for each of these variables. Each  $D$  variable has expected value  $q$  and variance  $q(1 - q)$ . If the two random variables are correlated, we require four probabilities, representing the outcomes 00, 01, 10 and 11, where, for example, 01 represents bond 1 being 0 and bond 2 being 1. If we assign these probabilities as

$$\begin{aligned}q_{00} &= (1 - q)^2 + \rho q(1 - q) \\q_{01} &= q(1 - q)(1 - \rho) \\q_{10} &= q(1 - q)(1 - \rho) \\q_{11} &= q^2 + \rho q(1 - q),\end{aligned}$$

then each variable will have probability of default  $q$ , as before. For example, the probability of  $D_1$  is  $q_{10} + q_{11} = q$ . The two  $D$ 's will now be correlated with a correlation coefficient of  $\rho$ .<sup>8</sup>

The corresponding price and other results for these new probabilities are shown in Table 17.2 for the case  $\rho = 0.7$ .

The main effect of positive correlation is that the probability of a single default decreases, because the likelihood that two will default increases.

In principle, a similar structure could be applied to a CDO with many instruments, but the evaluation would normally require simulation.

<sup>8</sup> To see this, we calculate  $\text{cov}(D_1, D_2) = \overline{D_1 D_2} - \overline{D}_1 \overline{D}_2 = q_{11} - q^2 = \rho q(1 - q)$ . Also,  $\text{var}(D_1) = \text{var}(D_2) = q(1 - q)$ . Hence,  $\text{cov}(D_1, D_2)/[\sigma(D_1)\sigma(D_2)] = \rho$ .

**TABLE 17.2**  
**RESULT WHEN DEFAULTS ARE CORRELATED WITH A CORRELATION COEFFICIENT OF 0.7**

Defaults	Prob	Total payout	Amount to A	Amount to B
0	0.873	200	120	80
1	0.054	130	120	10
2	0.073	60	60	0
		One	A	B
Mean payout		93.00	115.620	70.380
Price		88.464	109.981	66.948
Promised yield		.123	.087	.178
Default prob		.01	.073	.127
Sigma*		21.00	15.608	25.284

## 17.12 Summary

Credit derivatives are designed to ameliorate the risk associated with an entity's failure to meet its financial obligations, generally referred to as default. The standard credit derivative is the credit default swap (CDS), which is structured largely like an insurance policy.

The earliest modern approach to evaluation of credit risk is the Merton model, which treats a bank loan to a company as including a put option granted to the company. Value (for either the bank or the company) can thus be computed using standard option pricing mechanics under the assumption that the firm's value is governed by geometric Brownian motion. This idea has been extended to the use of other stochastic processes to model firm value, and the general method provides the basis of some commercial rating methodologies. This class of methods that are based on the trajectory of a firm's value are termed **structural models**.

By contrast, the **intensity (reduced-form)** method assumes that there is an underlying parameter that governs the probability of default over any interval of time. This intensity parameter may vary with time and, in fact, may be stochastic. Its behavior is modeled much like that of a short rate of interest; indeed, in essence the intensity acts like a spread factor that augments the standard short rate of interest. A reduced-form model is characteristic of a given credit risk entity and is typically calibrated by matching the spreads predicted by the model to the spreads of other instruments of the same (or similar) entity. Specific analytic models for the intensity process have been postulated and in some cases lead, if not easily tractable, at least to computer-implementable solutions for survival probabilities and other variables. Frequently, Monte Carlo simulation is used to solve for these quantities. The basic approach requires the generation of random default times, which can be consuming of computational time. A modified version that keeps track only of the expected value

of the impact of defaults is more efficient. Another approach to solution is to use lattice models, which offer great flexibility and have the advantage of being capable of determining optimal decisions, as in the case of an option on a CDS.

Although there are several viable approaches to analyzing the credit risk of a single entity, it is much more difficult to treat bundles of assets that are likely to be correlated in some fashion. One approach is to use a copula, which preserves the individual marginal probabilities of survival and embeds these in a structure that has interdependence.

## Exercises

- (Default probability) For the HiTech example (17.1), find the risk-neutral probability of default.
- (End default) For the HiTech bond of Example 17.4, suppose that default is recognized only at maturity and that no restitution is made, that is, the default is ignored. What is the value of that bond?
- (Bigger) Suppose that the firm of Example 17.6 has a value of \$200,000 instead of \$100,000. Perhaps this makes the bond more secure. What, in fact, is the value of the bond in this case?
- (Transitions) Suppose there are only 3 rating categories: A, B, D. Their one-year transition probabilities are given by the following matrix.

	Ending		
	Initial	A	B
A	.8	.2	0
B	.2	.7	.1

Construct the transition probability matrix for two years.

- (Average time overall) Using the density function of the stopping time probability for a fixed  $\lambda$ , find the average time to the first event over the entire interval  $[0, \infty)$ .
- (Fixed  $\lambda$ ) Consider a 3-year 10% coupon bond. The underlying short rate of interest follows a lattice with initial value of  $R = 1.15$  and then has an up factor of 1.02, a down factor of .99, and risk-neutral probabilities of .5.
  - Find the value of this bond with no risk of default.
  - Given that the risk-neutral intensity is constant such that the default probability is .9 each year, find the value of this bond.
  - What is the spread between the two bonds?
- (Sliding recovery) Consider a situation where  $r$  and  $\lambda$  are constant. A zero-coupon bond has face value  $F$  and maturity  $T$ . In the case of default at  $t$ , there is partial recovery equal to  $e^{-r(T-t)}F$ . Find the value of this bond.
- (Forward default rate) Let  $p(t)$  be the probability of survival from 0 to  $t$ . The probability of surviving to  $s$  given survival to  $t$  is then

$$p(s|t) = \frac{p(s)}{p(t)}.$$

Let

$$f(t) = -\frac{p'(t)}{p(t)}.$$

Show that

$$p(s|t) = e^{\int_t^s f(u)du}.$$

- 9.** (Default simulation) Refer to Section 17.8. Let  $q(t)$  be the survival probability and let  $q^{-1}$  be its inverse function. Also, let  $U$  be a uniform random variable on  $[0, 1]$ . For each realization  $u$ , let  $\tau$  be chosen such that  $q(\tau) = u$ . In this chain of relations, state whether each \* should be  $\leq$ ,  $=$ , or  $\geq$ :  $P[\tau * t] = P[q^{-1}(U) * t] = P[U * q(t)]$ . Conclude that since  $U$  is uniform,  $P[\tau \geq t] = q(t)$ .
- 10.** (10-year bond) Consider a 10-year zero-coupon bond with face value \$100. The interest rate is fixed at 5%. The credit spread for the bond is estimated to be 1% (except in part (a)). Calculate:

- (a) The bond value if there is no possibility of default.
- (b) The probability of default.
- (c) The value if there is no recovery.
- (d) The value if there is 50% recovery at maturity.
- (e) The value of \$100 that is paid at default.

[Hint: For (a),  $e^{-rT}F = \$6.65$ .]

- 11.** (Difficult estimation) Suppose we wish to estimate the probability of a rare event (such as a default probability). Let the random variable  $X$  be equal to 1 if the event occurs and to zero otherwise. Then  $p = E[X]$ . The standard Monte Carlo method takes  $n$  samples  $x_i$  and forms the average

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i.$$

- (a) What is the variance of  $\hat{p}$ ?
  - (b) What is the standard deviation of the absolute error  $\varepsilon_{ab} = \hat{p} - p$ ?
  - (c) Define the standard deviation of the relative error as  $\varepsilon_{rel} = \varepsilon_{ab}/p$ . Suppose  $p = 1\%$ . About how many samples must be taken to reduce the relative error to 1%.
- 12.** (TRORS difficulty) Consider a total return swap of a coupon bond versus a fixed-rate payment. Discuss whether the payoff is path dependent.
- 13.** (CDO variation) In Example 17.11, assume that the promised amounts to A and B are \$110 and \$90, respectively. Develop the new table of results. What are the prices of A and B?

## References

The Merton model of bankruptcy was first suggested in [1] and later elaborated in [2]. These two articles led to the structural approach to credit risk. The idea of using a crossing barrier was proposed in [3]. The nice lattice method for evaluation of coupon bonds, together with extensions, is in [4]. An empirical study of structural models is [5]. An introductory discussion of credit risk, including the approach to CDO's of Section 17.11 is [6]. An important review of credit risk is [7]. Three important and authoritative texts on credit risk are [8], [9], and [10]. This chapter was prepared by reference to those excellent texts and the overviews [11] and [12]. Also see [13] for methods to measure corporate default risk.

1. Black, F., and M. Scholes (1973), "The Pricing of Options and Corporate Liabilities," *Journal of Political Economy*, **81**, 637–654.
2. Merton, R. C. (1974), "On the Pricing of Corporate Debt: The Risk Structure of Interest Rates," *Journal of Finance*, **29**, 449–470.
3. Black, F., and J. Cox (1976), "Valuing Corporate Securities: Some Effects of Bond Indenture Provisions," *Journal of Finance*, **31**, 351–367.
4. Broadie, M., and Özgür Kaya (2007), "A Binomial Lattice Method for Pricing Corporate Debt and Modeling Chapter 11 Proceedings," *Journal of Financial and Quantitative Analysis*, **42**, no. 2, 279–312.
5. Eom, Y. H. , J. Helwege, and Jing-Zhi Huang (2004), "Structural Models of Corporate Bond Pricing: An Empirical Analysis," *The Review of Financial Studies*, **17**, no. 2, 499–544.
6. McDonald, Robert L. (2009), *Fundamentals of Derivatives Markets*, Prentice Hall, Boston.
7. Jarrow, R. A. (2009), "Credit Risk Models," *Annual Review of Financial Economics*, **1**: 37–68.
8. Duffie, Darrell and Kenneth J. Singleton (2003), *Credit Risk: Pricing, Measurement, and Management*, Princeton University Press, Princeton, NJ.
9. Schönbucher, Philipp J. (2003), *Credit Derivatives Pricing Models: Models, Pricing, and Implementation*, Wiley, Chichester, England.
10. Lando, David (2004), *Credit Risk Modeling*, Princeton University Press, Princeton, NJ.
11. Giesecke, Kay (April 2009), "An Overview of Credit Derivatives," Department of Management Science & Engineering, Also: (2009), *Jahresbericht der Deutschen Mathematiker-Vereinigung*, **111**.
12. Giesecke, Kay (October 2002), "Credit Risk Modeling and Valuation: An Introduction," Cornell University, Ithaca or, NY. Also abridged in D. Shimko (ed.) (2004), *Credit Risk: Models and Management*, Vol. 2, Riskbooks, London.
13. Duffie, Darrell (2011), *Measuring Corporate Default Risk*, Oxford University Press, New York.

## PART IV

# GENERAL CASH FLOW STREAMS





# 18

## OPTIMAL PORTFOLIO GROWTH

**C**onclusions about multiperiod investment situations are not mere variations of single-period conclusions—rather they often *reverse* those earlier conclusions. This makes the subject exciting, both intellectually and in practice. Once the subtleties of multiperiod investment are understood, the reward in terms of enhanced investment performance can be substantial.

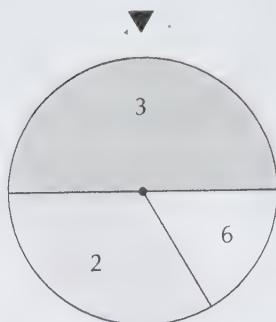
Fortunately the concepts and methods of analysis for multiperiod situations build on those of earlier chapters. Internal rate of return, present value, the comparison principle, portfolio design, and lattice and tree valuation all have natural extensions to general situations. But conclusions such as volatility is “bad” or diversification is “good” are no longer universal truths. The story is much more interesting.

This chapter begins the story by extending the elementary concept of internal rate of return, showing how to design portfolios that have maximal growth. The next chapter extends present value analysis.

### 18.1 The Investment Wheel

Understanding portfolio growth requires that one adopt a long-term viewpoint. To highlight the importance of such a viewpoint, consider the investment wheel shown in Figure 18.1. You are able to place a bet on any of the three sectors of the wheel. In fact, you may invest different amounts on each of the sectors independently. The numbers in the sectors denote the winnings for that sector after the wheel is spun.

**FIGURE 18.1 The investment wheel.** The numbers shown are the payoffs for a one-unit investment on that sector. The wheel is favorable and can be expected to cause capital to grow if investments are properly managed.



For example, if the wheel stops with the pointer at the top sector after a spin, you will receive \$3 for every \$1 you invested on that sector (which means a net profit of \$2).

The top sector is very attractive, paying 3 to 1, even though the area of that sector is a full one-half of the entire wheel. A \$1 bet (or investment) will return either \$0 or \$3, each with a probability of one-half. The expected gain is therefore  $\frac{1}{2} \times \$3 + \frac{1}{2} \times \$0 - \$1 = \$0.50$ . This is quite favorable.

The lower left sector, on the other hand, has unfavorable odds, since it pays only 2 to 1 for an area that is only one-third of the total. A bit better is the lower right segment, which pays even odds, since it pays 6 to 1 and is one-sixth of the area.

Suppose now that you start with \$100 and have the opportunity to bet part or all of your money repeatedly, reinvesting your winnings on successive spins of the wheel. Because of the favorable top segment, you can make your capital grow over the long run through judicious investment. The question is, just what constitutes judicious investment?

Based on the odds we calculated, it seems appropriate to concentrate your attention (and your capital) on the top sector. One strategy would be to invest all of your money on that sector. Indeed, this strategy is the one that produces the highest single-period expected return. An investment of \$100 is expected to gain an additional \$50 on the very first spin. The problem is that you go broke half of the time and cannot continue with other spins. Even if you win and continue with this strategy, you will again face the risk of ruin at the next spin. Most people find this strategy too risky when given the opportunity to play repeatedly.

A second, more conservative, strategy would be to invest, say, one-half of your money on the top sector each spin, holding back the other half. That way if an unfavorable outcome occurs, you are not out of the game entirely.<sup>1</sup> But it is not clear that this is the best that can be done.

<sup>1</sup> This wheel investment problem actually makes a good game for a group, using play money or keeping records. Actual play forces people to think exactly how they wish to invest. The main point is that investment for the long run is not the same as investment for a single spin.

## Analysis of the Wheel

To begin a systematic search for a good strategy, let us limit our investigation to **fixed-proportions** strategies. These are strategies that prescribe proportions to each sector of the wheel, these proportions being used to apportion current wealth among the sectors as bets at each spin. Let us number the sectors 1, 2, and 3, corresponding to top, left, and right, respectively. A general fixed-proportions strategy for the wheel is then described by a set of three numbers  $(\alpha_1, \alpha_2, \alpha_3)$ , where each  $\alpha_i \geq 0$ ,  $i = 1, 2, 3$ , and where  $\alpha_1 + \alpha_2 + \alpha_3 \leq 1$ . The  $\alpha_i$ 's correspond to the proportions bet on the different sectors. The remaining  $1 - \alpha_1 - \alpha_2 - \alpha_3$  is held in reserve. As an example, the strategy mentioned earlier of investing one-half of your capital in the top segment each time is  $(\frac{1}{2}, 0, 0)$ .

Each fixed-proportions strategy leads to a series of multiplicative factors that govern the growth of capital. For example, suppose you bet \$100 using the  $(\frac{1}{2}, 0, 0)$  strategy. For one spin there are two possibilities: (1) with probability one-half you obtain a favorable outcome and end up with  $\$50 + 3 \times \$50 = \$200$ ; and (2) with probability one-half you obtain an unfavorable outcome and end up with just \$50. In general, with this strategy your money will be either doubled or halved at each spin, each possibility occurring with probability one-half. The multiplicative factors for one spin are thus 2 and  $\frac{1}{2}$ , each with probability one-half. After a long series of investments following this strategy, your initial \$100 will be multiplied by an overall multiple that might be of the form  $(\frac{1}{2})(\frac{1}{2})(2)(\frac{1}{2})(2)(2) \dots (2)(\frac{1}{2})$ , with about an equal number of 2's and  $\frac{1}{2}$ 's. Hence the overall factor is likely to be about 1. This means that during the course of many spins, your capital will tend to fluctuate up and down, but is unlikely to grow appreciably.

An alternative strategy is to bet one-fourth of your money on the top sector, corresponding to the strategy  $(\frac{1}{4}, 0, 0)$ . If that top sector is the outcome of a spin, your money will be multiplied by  $1 - \frac{1}{4} + \frac{3}{4} = \frac{3}{2}$ . If that sector is not the outcome, your money will be multiplied by  $1 - \frac{1}{4} = \frac{3}{4}$ . On average, two spins provide a factor of  $(\frac{3}{2})(\frac{3}{4}) = \frac{9}{8}$ . Hence each single spin provides, on average, a factor of  $\sqrt{\frac{9}{8}} = 1.06066$ . With this strategy your money will grow, on average, by over 6% each turn. (Exercise 1 shows that this strategy is, in a limited sense, optimal.)

## 18.2 The Log Utility Approach to Growth

The investment wheel is representative of a large and important class of investment situations where a particular strategy leads to a random growth process. This class includes investment in common stocks, as shown later in this section. A general formulation is that if  $X_k$  represents capital after the  $k$ th trial, then

$$X_k = R_k X_{k-1} \quad (18.1)$$

for  $k = 1, 2, \dots$ . In this equation  $R_k$  is a random return variable. We assume that it is a **stationary independent** process, where all  $R_k$ 's have identical probability distributions and are mutually independent.

The investment wheel with the strategy of investing one-half of one's capital on the top segment corresponds to this model with  $R_k$ 's that take on either of the two values 2.0 or .50, each with probability of one-half. The  $R_k$  variables all have the same probability density and are independent of one another (that is, other outcomes do not influence the present outcome).

In the general capital growth process, the capital at the end of  $n$  trials is

$$X_n = R_n R_{n-1} \cdots R_2 R_1 X_0.$$

Taking the logarithm of both sides gives

$$\ln X_n = \ln X_0 + \sum_{k=1}^n \ln R_k.$$

A little more manipulation produces

$$\ln \left( \frac{X_n}{X_0} \right)^{1/n} = \frac{1}{n} \sum_{k=1}^n \ln R_k. \quad (18.2)$$

Consider the right-hand side of equation (18.2) as  $n \rightarrow \infty$ . The variables  $\ln R_k$  are each random variables that are independent and have identical probability distributions. **The law of large numbers**<sup>2</sup> therefore states that

$$\frac{1}{n} \sum_{k=1}^n \ln R_k \rightarrow E(\ln R_1).$$

(We can use  $E(\ln R_1)$  in this expression since the expected value is the same for all  $k$ .)

We define  $m = E(\ln R_1)$ . Then we have, from equation (18.2),

$$\ln \left( \frac{X_n}{X_0} \right)^{1/n} \rightarrow m.$$

This is the fundamental result that we now highlight:

**Logarithmic performance** *If  $X_1, X_2, \dots$  is the random sequence of capital values generated by the process*

$$X_k = R_k X_{k-1},$$

*then*

$$\ln \left( \frac{X_n}{X_0} \right)^{1/n} \rightarrow m \quad (18.3)$$

---

<sup>2</sup> The law of large numbers states that if  $Y_1, Y_2, \dots$  are independent random variables with identical distributions then  $(1/n) \sum_{k=1}^n Y_k \rightarrow E(Y_k)$ . A simple example is that of flipping a coin and assigning  $Y_k = +1$  if heads occurs on the  $k$ th trial and  $-1$  if tails occurs. The average of the numbers tends to zero.

as  $n \rightarrow \infty$ , where

$$m = E(\ln R_1). \quad (18.4)$$

Taking the antilogarithm of both sides of (18.3) gives

$$\left( \frac{X_n}{X_0} \right)^{1/n} \rightarrow e^m.$$

Then, formally (although it is not quite legitimate to do so), we raise both sides to the power of  $n$ , and we find

$$X_n \rightarrow X_0 e^{mn}.$$

In other words, for large  $n$  the capital grows (roughly) exponentially with  $n$  at a rate  $m$ .

The foregoing analysis reveals the importance of the number  $m$  defined by (18.4). It governs the rate of growth of the investment over a long period of repeated trials. It seems appropriate therefore to select the strategy that leads to the largest value of  $m$ .

## Log Utility Form

Note that if we add the constant  $\ln X_0$  to (18.4) we find

$$m + \ln X_0 = E(\ln R_1) + \ln X_0 = E(\ln R_1 X_0) = E(\ln X_1).$$

Hence if we define the special utility function  $U(X) = \ln X$ , the problem of maximizing the growth rate  $m$  is equivalent to maximizing the expected utility  $E[U(X_1)]$  and using this same strategy in every trial. In other words, by using the logarithm as a utility function, we can treat the problem as if it were a single-period problem. We find the optimal growth strategy by finding the best thing to do on the first trial, using the expected logarithm as our criterion. This single-step view guarantees the maximum growth rate in the long run.

## Examples

Many important and interesting situations fit the framework presented in this section.

**Example 18.1 (The Kelly rule of betting)** Suppose that you have the opportunity to invest in a prospect that will either double your investment or return nothing. The probability of the favorable outcome is  $p$ . Suppose that you have an initial capital of  $X_0$  and you can repeat this investment many times. How much should you invest each time?

This situation closely resembles the game of blackjack, played by a player who mentally keeps track of the cards played. By adjusting the playing strategy to account for the composition of the remaining deck, such a player may have, on average, about a 50.75% chance of winning a hand; that is,  $p = .5075$ . The player must decide how much to bet in such a situation.

Let  $\alpha$  be the proportion of capital invested (or bet) during one play. The player wishes to find the best value of  $\alpha$ . If the player wins, his or her capital will grow by the factor  $1 - \alpha + 2\alpha = 1 + \alpha$ . If he or she loses, the factor is  $1 - \alpha$ . Hence to find the log-optimal value of  $\alpha$ , we maximize

$$m = p \ln(1 + \alpha) + (1 - p) \ln(1 - \alpha).$$

Setting the derivative with respect to  $\alpha$  equal to zero, we have

$$\frac{p}{1 + \alpha} - \frac{1 - p}{1 - \alpha} = 0.$$

This gives the equation

$$p(1 - \alpha) - (1 - p)(1 + \alpha) = 0,$$

or  $\alpha = 2p - 1$ .<sup>3</sup> Hence in the blackjack example, a player should bet 1.5% of the total capital on each hand when  $p = .5075$ . Professional blackjack players actually do use this rule or a modification of it.

Blackjack may seem to offer an easy living! The growth rate of the Kelly rule strategy is

$$m = p \ln 2p + (1 - p) \ln(2 - 2p) = p \ln p + (1 - p) \ln(1 - p) + \ln 2.$$

For the case where  $p = .5075$ , this gives  $e^m \approx 1.0001125$ , which is a .01125% gain. To double your capital you must expect to play  $72/.01125 = 6,440$  hands (remember the 72 rule of Chapter 2). This requires about 80 hours of play, which realistically requires about 1 month of activity. But there are many obstacles in the path of such a profession.

**Example 18.2 (Volatility pumping)** Suppose there are two assets available for investment. One is a stock that in each period either doubles or reduces by one-half, each with a probability of 50%. The other just retains value—like putting money under the mattress. Neither of these investments is very exciting. An investment left in the stock will have a value that fluctuates a lot but has no overall growth rate. The other clearly has no growth rate. Nevertheless, by using these two investments in combination, growth can be achieved.

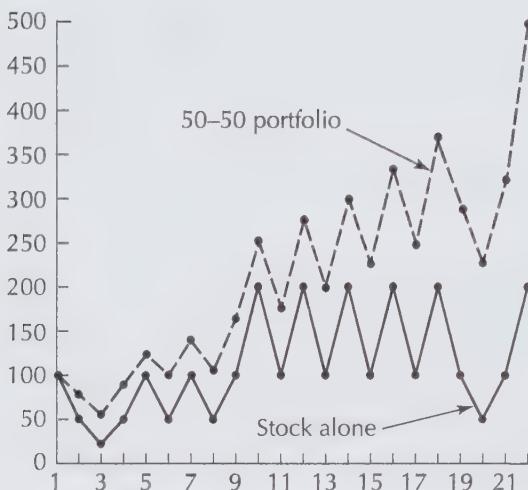
To see how, suppose that we invest one-half of our capital in each asset each period. Thus we **rebalance** at the beginning of each period by being sure that one-half of our capital is in each asset. Under a favorable performance, our capital will grow by the factor  $\frac{1}{2} + \frac{1}{2} \times 2 = \frac{1}{2} + 1$ . Under an unfavorable performance, the factor will be  $\frac{1}{2} + \frac{1}{2} \times \frac{1}{2} = \frac{1}{2} + \frac{1}{4}$ . Hence the expected growth rate of this strategy is

$$m = \frac{1}{2} \ln(\frac{1}{2} + 1) + \frac{1}{2} \ln(\frac{1}{2} + \frac{1}{4}) \approx .059.$$

Therefore  $e^m = 1.0607$ , and the gain on the portfolio is about 6% per period.

---

<sup>3</sup> The answer implicitly assumes  $p > .5$ . If  $p \leq .5$ , the optimal  $\alpha$  is  $\alpha = 0$ .



**FIGURE 18.2 Mixture of two assets.** Two mediocre stocks can be combined to give enhanced growth.

Figure 18.2 shows one simulation run of the performance of the 50–50 mix of the two assets versus the stock itself. The mixture portfolio outperforms the stock.

The gain is achieved by using the volatility of the stock in a **pumping** action. If the stock goes up in a certain period, some of the proceeds are put aside. If on the other hand the stock goes down, additional capital is invested in it. Capital is pumped back and forth between the two assets in order to achieve growth greater than can be achieved by either alone.

Note also that this strategy automatically, on average, follows the dictum of “buy low and sell high” by the process of rebalancing. In essence, that is why it produces growth.

**Example 18.3 (Pumping two stocks)** Let us modify Example 18.2 by assuming that both assets have the property of either doubling or halving in value each period with probability one-half. Each asset moves independently of the other. Again we invest one-half of our capital in each asset, rebalancing at each period. We find immediately that

$$m = \frac{1}{4} \ln 2 + \frac{1}{2} \ln \frac{5}{4} + \frac{1}{4} \ln \frac{1}{2} = \frac{1}{2} \ln \frac{5}{4} = .1116.$$

Hence  $e^m = \sqrt{\frac{5}{4}} = 1.118$ , which corresponds to an 11.8% growth rate each period. The pumping action is greatly enhanced over that of the previous example. Pumping between two volatile assets leads to large growth rates.

**Example 18.4 (Large stock portfolios)** Suppose that there are  $n$  stocks that have returns  $R_i$ ,  $i = 1, 2, 3, \dots, n$ , for any one period (of, say, a week). These returns are random, but they have the same probability distribution each period. The returns

of different stocks may be correlated, but the returns of different periods are not correlated. We form a portfolio of these stocks by assigning weights  $w_1, w_2, w_3, \dots, w_n$  with  $w_i \geq 0$  for each  $i$  and  $\sum_{i=1}^n w_i = 1$ . The overall return on the portfolio is  $R = \sum_{i=1}^n w_i R_i$ . To obtain the maximum possible growth of this portfolio, we select the weights so as to maximize  $m = E(\ln R)$ . If we do this, the portfolio can be expected to grow roughly, on average, according to  $e^{mk}$ , where  $k$  is the number of periods.

We shall study this example in greater detail later in this chapter.

**Example 18.5 (The investment wheel)** Let us compute the full optimal strategy for the investment wheel allowing for the possibility of investing on all sectors. For a strategy<sup>4</sup>  $(\alpha_1, \alpha_2, \alpha_3)$  we find the results as follows:

1. If 1 occurs,  $R = 1 + 2\alpha_1 - \alpha_2 - \alpha_3$ .
2. If 2 occurs,  $R = 1 - \alpha_1 + \alpha_2 - \alpha_3$ .
3. If 3 occurs,  $R = 1 - \alpha_1 - \alpha_2 + 5\alpha_3$ .

To maximize the expected logarithm of this return structure, we maximize

$$m = \frac{1}{2} \ln(1 + 2\alpha_1 - \alpha_2 - \alpha_3) + \frac{1}{3} \ln(1 - \alpha_1 + \alpha_2 - \alpha_3) + \frac{1}{6} \ln(1 - \alpha_1 - \alpha_2 + 5\alpha_3).$$

If we assume that the solution has  $\alpha_i > 0$  for each  $i = 1, 2, 3$ , we can find the solution by setting the derivatives with respect to each  $\alpha_i$  equal to zero. This gives the equations

$$\begin{aligned} \frac{2}{2(1 + 2\alpha_1 - \alpha_2 - \alpha_3)} - \frac{1}{3(1 - \alpha_1 + \alpha_2 - \alpha_3)} - \frac{1}{6(1 - \alpha_1 - \alpha_2 + 5\alpha_3)} &= 0 \\ -\frac{1}{2(1 + 2\alpha_1 - \alpha_2 - \alpha_3)} + \frac{1}{3(1 - \alpha_1 + \alpha_2 - \alpha_3)} - \frac{1}{6(1 - \alpha_1 - \alpha_2 + 5\alpha_3)} &= 0 \\ -\frac{1}{2(1 + 2\alpha_1 - \alpha_2 - \alpha_3)} - \frac{1}{3(1 - \alpha_1 + \alpha_2 - \alpha_3)} + \frac{5}{6(1 - \alpha_1 - \alpha_2 + 5\alpha_3)} &= 0. \end{aligned}$$

General equations of this form are difficult to solve analytically. However, in this case a solution is  $\alpha_1 = \frac{1}{2}$ ,  $\alpha_2 = \frac{1}{3}$ , and  $\alpha_3 = \frac{1}{6}$ , which can be checked easily. (For a generalization of this problem and its solution see Exercise 4.) This means that one should invest in every sector of the wheel, and the proportions bet are equal to the probabilities of occurrence.

Substitution of this optimal strategy in the original objective of expected logarithm gives

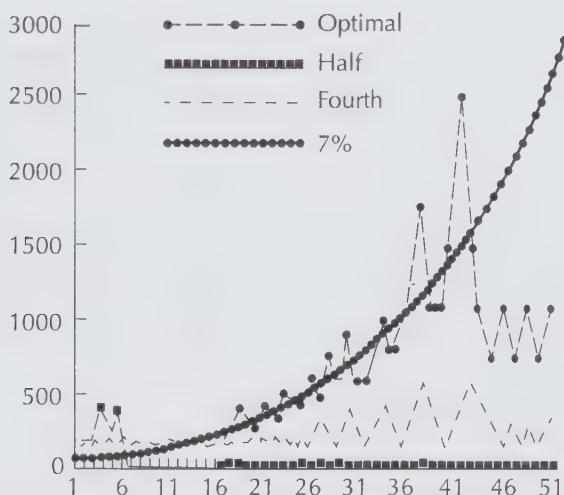
$$m = \frac{1}{2} \ln \frac{3}{2} + \frac{1}{3} \ln \frac{2}{3} + \frac{1}{6} \ln 1 = \frac{1}{6} \ln \frac{3}{2}.$$

We then find that

$$e^m \approx 1.06991.$$

Hence the optimal solution achieves a growth rate of about 7%, which compares with the approximately 6% achieved by the strategy of investing one-fourth on the top segment and nothing on the other two.

<sup>4</sup> Recall that 1, 2, 3 correspond to the top, left, and right, with payoffs 3, 2, and 6, respectively.



**FIGURE 18.3** **Wheel simulation.** Under the optimal strategy, the wheel provides a growth rate of nearly 7%.

The results of one simulation of 50 trials of the investment wheel are shown in Figure 18.3. The figure shows the results for three strategies: the optimal strategy, the simplified strategy of betting one-fourth on the top segment, and the poor strategy of investing one-half on the top segment. Also shown is a curve representing a 7% growth rate. The simulation has a great deal of volatility, and other runs may look quite different from this one. The long-term effect shows up when there are hundreds of trials, as there would be, for example, in the yearly result of daily stock market investments.

Notice that the optimal strategy requires an investment on the unfavorable sector 2, which pays only 2 to 1. This investment serves as a hedge for the other sectors—it wins precisely when the others do not.<sup>5</sup> It is like fire insurance on your home, paying when other things go wrong.

## 18.3 Properties of the Log-Optimal Strategy\*

Although the log-optimal strategy maximizes the expected growth rate, the short run growth rate may differ. We can, however, make some definite statements about the log-optimal strategy that are quite impressive.

<sup>5</sup> The equations defining the optimal solution are actually degenerate for this problem. There is a whole family of optimal solutions, all giving the same value for  $m$ . An alternate solution is  $\alpha_1 = \frac{5}{18}$ ,  $\alpha_2 = 0$ ,  $\alpha_3 = \frac{1}{18}$ . In this solution nothing is invested on the unfavorable sector.

Suppose two people start with the same initial capital level  $X_0$ . Suppose further that person A invests using the log-optimal strategy and person B uses some other strategy (with a lower value of  $m$ ). Denote the resulting capital streams by  $X_k^A$  and  $X_k^B$ , respectively, for the periods  $k = 1, 2, \dots$ . Then it can be shown that

$$E(X_k^B / X_k^A) \leq 1 \quad \text{for all } k.$$

This says that the ratio of the capital associated with alternative strategy B to the capital associated with the optimal strategy A is expected to be less than 1 at every stage. This property argues in favor of using the log-optimal strategy, and many people are indeed persuaded that this is the strategy they should adopt.

## 18.4 Alternative Approaches\*

The log-optimal strategy is not necessarily the best strategy to use in repetitive investment situations, but it is a good benchmark to keep in mind when considering alternatives. We mention some possible alternatives in this section.

### Other Utility

One alternative is to use the standard framework of maximizing expected utility (as in the first part of Chapter 11). If there will be exactly  $K$  repetitions, we can define a utility function  $U$  for wealth at the end of period  $K$  and, accordingly, seek to maximize  $E[U(X_K)]$ .

The use of  $U(X_K) = \ln X_K$  is one special case. In fact, because of a special recursive property, maximization of  $E(\ln X_K)$  with respect to a fixed strategy is exactly equivalent to the log-optimal strategy of maximizing  $E(\ln X_1)$ . This follows from

$$E(\ln X_K) = E[\ln(R_K R_{K-1} \cdots R_1 X_0)] = \ln X_0 + E(\ln R_1) + \sum_{k=2}^K E(\ln R_k).$$

Maximization of the left side is equivalent to maximization of  $E(\ln R_1)$  since all  $R_k$ 's are identical random variables. This in turn is equivalent to maximization of  $E(\ln X_1) = \ln X_0 + E(\ln R_1)$ . Hence the choice of  $U(X_K) = \ln X_K$  leads us once again to the log-optimal strategy.

One interesting class of utility functions is the class of power functions  $U(X) = (1/\gamma)X^\gamma$  for  $\gamma \leq 1$ . This class includes the logarithm [since  $\lim_{\gamma \rightarrow 0} \{(1/\gamma)X^\gamma - 1/\gamma\} = \ln X$ ]; and it includes the linear utility  $U(X) = X$ .

This class of functions has the same recursive property as the log utility; that is, the structure is preserved from period to period. This is seen from

$$\begin{aligned} E[U(X_K)] &= \frac{1}{\gamma} E[(R_K R_{K-1} \cdots R_1 X_0)^\gamma] = \frac{1}{\gamma} E(R_K^\gamma R_{K-1}^\gamma \cdots R_1^\gamma) X_0^\gamma \\ &= \frac{1}{\gamma} E(R_K^\gamma) E(R_{K-1}^\gamma) \cdots E(R_1^\gamma) X_0^\gamma \end{aligned}$$

where the last equality follows from the fact that the expected value of a product of independent random variables is equal to the product of the expected values. Hence to maximize  $E[U(X_K)]$  with a fixed-proportions strategy it is only necessary to maximize  $E[(R_1 X_0)^\gamma]$ , so again to maximize  $E[U(X_K)]$  one need only maximize  $E[U(X_1)]$ .

If  $\gamma > 0$ , the power utility function is quite aggressive. The extreme case of  $\gamma = 1$ , corresponding to  $U(X) = X$  (leading to the expected-value criterion), was considered earlier when discussing the investment wheel. We found that the strategy that maximizes the expected value bets all capital on the most favorable sector—a strategy prone to early bankruptcy. Indeed, bankruptcy is likely for any  $\gamma$  with  $1 \geq \gamma > 0$ . For example, suppose  $\gamma = \frac{1}{2}$ . Consider two opportunities: (a) capital will double with a probability of .90 or it will go to zero with probability .10, and (b) capital will increase by 25% with certainty. Since  $.9 \times \sqrt{2} > \sqrt{1.25}$ , opportunity (a) is preferred to (b) with a square root utility. However, in a long sequence of repeated trials, an investor following opportunity (a) is virtually certain to go bankrupt. Most people prefer (b) when they understand that many trials will be played. A similar argument applies to any  $\gamma$  in the range  $1 \geq \gamma > 0$ .

It is more conservative to use  $\gamma < 0$ . However, many people find this to be *too* conservative. For example, suppose that  $\gamma = -\frac{1}{2}$ . Again consider two opportunities: (a) capital quadruples in value with certainty, and (b) with probability .5 capital remains constant and with probability .5 capital is multiplied by 10 million (or any finite number). Since  $-4^{-1/2} > -.5 - .5(10,000,000)^{-1/2}$ , an investor with the utility function  $V(X) = -X^{-1/2}$  will prefer (a). This is quite conservative. Again, similar arguments apply for any  $\gamma < 0$ , although they become less compelling if  $\gamma$  is close to zero.

Based on the preceding discussion, we conclude that if an investor uses a power utility function, it is likely that it will be one with  $\gamma < 0$ , but  $\gamma$  close to zero. Such a utility function is close to the logarithm. We can argue that similar (although less precise) results hold for any broad class of possible utility functions; that is, only those close to the logarithm will seem appropriate when the long-term consequences are examined. Therefore, although in principle an investor may choose any utility function, supposedly reflecting individual risk tolerance, a repetitive situation tends to hammer the utility into one that is close to the logarithm.

Most long-term investors do consider the volatility of portfolio growth as well as the growth rate itself. This leads to consideration of the variance of the logarithm of return as well as the expected value of the logarithm of return. Indeed, if investors take a long-term view, it can be shown that (under certain assumptions) these two values are the only values of importance. We state this formally as follows:

**Growth efficiency proposition** *An investor who considers only long-term performance will evaluate a portfolio on the basis of its logarithm of single-period return, using only the expected value and the variance of this quantity.*

This proposition interlaces well with the earlier discussion about power utility functions. We found that if the utility function  $U(X_1) = (1/\gamma)X^\gamma$  were chosen, it is

likely that  $\gamma < 0$  and  $\gamma \approx 0$ . We can then use the approximation

$$\frac{1}{\gamma}(X^\gamma - 1) \approx \ln X + \frac{1}{2}\gamma(\ln X)^2.$$

This shows that use of this utility function is close to using a weighted combination of the expected logarithm of return and the variance of that logarithm. In other words, the expected logarithm and its variance are the two quantities of interest.

In view of the growth efficiency proposition, it is natural to trace out an efficient frontier of  $m$  versus  $\sigma$  similar to that for the ordinary mean-variance efficient frontier but where  $m$  and  $\sigma$  are, respectively, the mean and standard deviation of the logarithm of return. We shall do this for stocks whose prices are described by continuous-time equations in the next section.

## 18.5 Continuous-Time Growth

Optimal portfolio growth can be applied with any rebalancing period—a year, a month, a week, or a day. In the limit of very short time periods we consider continuous rebalancing.

In fact, there is a compelling reason to consider the limiting situation: the resulting equations for optimal strategies turn out to be much simpler, and as a consequence it is much easier to compute optimal solutions. Hence even if rebalancing is to be carried out only, say, weekly, it is convenient to use the continuous-time formulation to do the calculations.

The continuous-time version also provides important insight. For example, it reveals very clearly how volatility pumping works.

## Dynamics of Several Stocks

We first extend the continuous-time model of stock dynamics presented in Chapter 13 to the case of several correlated stocks. This model will then be used in our analysis of stock portfolios.

Suppose there are  $n$  assets. The price  $p_i$  of the  $i$ th asset, for  $i = 1, 2, 3, \dots, n$ , is governed by a standard geometric Brownian motion equation

$$\frac{dp_i}{p_i} = \mu_i dt + dz_i$$

where  $z_i$  denotes a Wiener process, but with variance parameter  $\sigma_i^2$  rather than 1. This is equivalent to the standard model for a single stock. The new element here is that the assets are correlated through the Wiener process components. In particular,

$$\text{cov}(dz_i, dz_j) = E(dz_i dz_j) = \sigma_{ij} dt.$$

We define the **covariance matrix  $S$**  as that with components  $\sigma_{ij}$ , and we use the convention  $\sigma_i^2 = \sigma_{ii}$ . We usually assume that  $S$  is nonsingular.

From Chapter 13, each asset  $i$  has a lognormal distribution, and at time  $t$ ,

$$\mathbb{E} \left[ \ln \left( \frac{p_i(t)}{p_i(0)} \right) \right] = (\mu_i - \frac{1}{2} \sigma_i^2) t \equiv v_i t$$

and

$$\text{var} \left[ \ln \left( \frac{p_i(t)}{p_i(0)} \right) \right] = \sigma_i^2 t.$$

## Portfolio Dynamics

Now suppose that a portfolio of the  $n$  assets is constructed using the weights  $w_i$ ,  $i = 1, 2, \dots, n$ , with  $\sum_{i=1}^n w_i = 1$ . Let  $V$  be the value of the portfolio. Then because the instantaneous rate of return of the portfolio is equal to the weighted sum of the instantaneous rates of return of the individual assets, we have

$$\begin{aligned} \frac{dV}{V} &= \sum_{i=1}^n w_i \frac{dp_i}{p_i} \\ &= \sum_{i=1}^n w_i \mu_i dt + w_i dz_i. \end{aligned}$$

The variance of the stochastic term is

$$\mathbb{E} \left( \sum_{i=1}^n w_i dz_i \right)^2 = \mathbb{E} \left[ \left( \sum_{i=1}^n w_i dz_i \right) \left( \sum_{j=1}^n w_j dz_j \right) \right] = \sum_{i,j=1}^n w_i \sigma_{ij} w_j dt.$$

Hence the value  $V(t)$  is lognormal with

$$\mathbb{E} \left[ \ln \left( \frac{V(t)}{V(0)} \right) \right] = vt = \sum_{i=1}^n w_i \mu_i t - \frac{1}{2} \sum_{i,j} w_i \sigma_{ij} w_j t. \quad (18.5)$$

The variance of  $\ln[V(t)/V(0)]$  is

$$\sigma^2(t) = \sum_{i,j}^n w_i \sigma_{ij} w_j t.$$

Note that

$$\nu = \frac{1}{t} \mathbb{E} \left[ \frac{\ln V(t)}{V(0)} \right].$$

Hence  $\nu$  gives the growth rate of the portfolio—analogous to  $m$ , used in previous sections. We can control this growth rate by the choice of the weighting coefficients  $w_1, w_2, \dots, w_n$ .

## Implications for Growth

Equation (18.5) explains how volatility can be pumped to obtain increased growth. As a specific example, suppose that the  $n$  assets are uncorrelated and all have the same mean and variance. A typical asset therefore has its price governed by the process

$$\frac{dp_i}{p_i} = \mu dt + dz_i$$

where now each  $dz_i$  has variance  $\sigma^2 dt$ . The expected growth rate of each stock individually is  $v = \mu - \frac{1}{2}\sigma^2$ . Suppose now that the  $n$  stocks are each included in a portfolio with a weight of  $1/n$ . Then from equation (18.5) the expected growth rate of the portfolio is

$$v_{\text{port}} = \mu - \frac{1}{2n}\sigma^2.$$

Pumping reduces the magnitude of the  $-\frac{1}{2}\sigma^2$  correction term, thereby increasing the growth rate. In this example, the growth rate has increased over the  $v$  of a single stock by

$$v_{\text{port}} - v = \frac{1}{2} \left(1 - \frac{1}{n}\right)\sigma^2 = \frac{1}{2} \left(\frac{n-1}{n}\right)\sigma^2.$$

The pumping effect is obviously most dramatic when the original variance is high. After being convinced of this, you will likely begin to *enjoy* volatility, seeking it out for your investments rather than shunning it, as you may have after studying the single-period theory of Chapters 6 and 7. Volatility is *not* the same as risk. Volatility is opportunity.

**Example 18.6 (Volatility in action)** Suppose that a stock has an expected growth rate of 15% a year and a volatility (of its logarithm) of 20%. These are fairly typical values. This means that  $v = \mu - \frac{1}{2}\sigma^2 = .15$  and  $\sigma = .20$ . Hence  $\mu = .15 + .04/2 = .17$ . By combining 10 such stocks in equal proportions (and assuming they are uncorrelated) we obtain an overall growth rate improvement of  $(9/20) \times .04 = 1.8\%$ —nice, but not dramatic.

If instead the individual volatilities were 40%, the improvement in growth rate would be 7.2%, which is substantial. At volatilities of 60% the improvement would be 16.2%, which is truly impressive. Unfortunately it is hard to find 10 uncorrelated stocks with this level of volatility, so in practice one must settle for more modest gains.<sup>6</sup>

## The Portfolio of Maximum Growth Rate

We obtain the optimal growth portfolio by maximizing the growth rate  $v$ . Referring to equation (18.5) we see that this is accomplished by finding the weights  $w_1, w_2, \dots, w_n$

---

<sup>6</sup> Of course we must temper our enthusiasm with an accounting of the commissions associated with frequent trading.

that solve

$$\begin{aligned} & \text{maximize} \sum_{i=1}^n w_i \mu_i - \frac{1}{2} \sum_{i,j=1}^n w_i \sigma_{ij} w_j \\ & \text{subject to} \sum_{i=1}^n w_i = 1. \end{aligned}$$

We solve this problem in the next section.

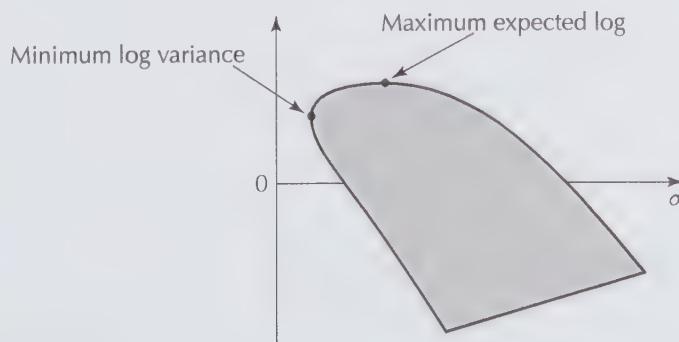
## 18.6 The Feasible Region

Paralleling the familiar Markowitz concept, portfolios can be plotted on a two-dimensional diagram of  $\nu$  versus  $\sigma$ . The region mapped out by all possible portfolios defines the feasible region. This is depicted in Figure 18.4.

There is, however, an important qualitative difference between the general shape of this region and the Markowitz region. The new region does not extend upward indefinitely, but instead there is a maximum value of  $\nu$ , corresponding to the growth rate of the log-optimal portfolio. There is also, as in the Markowitz case, a point of minimum  $\sigma$ . These points are indicated on the figure.

### The Efficient Frontier

Again, just as in the Markowitz framework, we define the **efficient frontier** of the feasible region to be the upper left-hand portion of the boundary. This frontier is efficient in the sense of growth as spelled out by the growth efficiency proposition of Section 18.4. In this case we can be quite specific and state that the efficient frontier is the portion of the boundary curve lying between the minimum-variance point and the log-optimal point.



**FIGURE 18.4 Feasible region.** The feasible region has a maximum expected log value and a minimum log variance value.

In fact, we obtain a strong version of the **two-fund theorem**. Any point on the efficient frontier can be achieved by a portfolio consisting of a mixture of the minimum-variance portfolio and the log-optimal portfolio. We now state this formally as a theorem. We also give a proof using vector-matrix notation. (The reader may safely skip the proof.)

**The two-fund theorem** *Any point on the efficient frontier can be achieved as a mixture of any two points on frontier. In particular the minimum-log-variance portfolio and the log-optimal portfolio can be used.*

**Proof:** Assume there are  $n$  securities. Let  $\mathbf{u} = (\mu_1, \mu_2, \dots, \mu_n)$ , and let  $\mathbf{w} = (w_1, w_2, \dots, w_n)$  be portfolio weights. If  $\mathbf{w}$  is efficient, it must solve the following problem for some  $s$ :

$$\begin{aligned} & \text{maximize } \mathbf{w}^T \mathbf{u} - \frac{1}{2} \mathbf{w}^T \mathbf{S} \mathbf{w} \\ & \text{subject to } \mathbf{w}^T \mathbf{1} = 1 \\ & \quad \mathbf{w}^T \mathbf{S} \mathbf{w} = s. \end{aligned}$$

By introducing Lagrange multipliers  $\lambda$  and  $\gamma/2$ , we form the Lagrangian

$$L = \mathbf{w}^T \mathbf{u} - \frac{1}{2} \mathbf{w}^T \mathbf{S} \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{1} - 1) - \frac{1}{2}\gamma(\mathbf{w}^T \mathbf{S} \mathbf{w} - s).$$

The first-order conditions are

$$\mathbf{u} - \mathbf{S} \mathbf{w} - \lambda \mathbf{1} - \gamma \mathbf{S} \mathbf{w} = \mathbf{0}.$$

Hence the solution has the form

$$\mathbf{w} = \frac{1}{1 + \gamma} \mathbf{S}^{-1}(\mathbf{u} - \lambda \mathbf{1}).$$

The constants  $\lambda$  and  $\gamma$  are determined so that the solution  $\mathbf{w}$  satisfies the two constraints of the original problem.

Setting  $\gamma = 0$  means that the second constraint is not active, and hence this solution corresponds to the log-optimal portfolio.

All solutions are linear combinations of the two vectors  $\mathbf{S}^{-1} \mathbf{u}$  and  $\mathbf{S}^{-1} \mathbf{1}$ , so any two such solutions can be used to generate all others. In particular, the log-optimal and the minimum-variance solutions can be used. ■

## Inclusion of a Risk-Free Asset

Suppose that there is a risk-free asset with constant interest rate  $r_f$ . This asset can be considered to be a bond whose price  $p_0(t)$  satisfies the equation

$$\frac{dp_0(t)}{p_0} = r_f dt.$$

Assuming that there is no other combination of assets that produces zero variance, the risk-free asset is on the efficient frontier. Indeed, it is the minimum-variance point. To find the entire efficient frontier it is therefore only necessary to find the log-optimal point, and we shall do that now.

The log-optimal portfolio is defined by a set of weights  $w_1, w_2, \dots, w_n$  for the risky assets and a weight  $w_0 = 1 - \sum_{j=1}^n w_j$  for the risk-free asset. The weights for the risky assets are chosen to maximize the overall growth rate; that is, to solve the problem

$$\max \left[ \left( 1 - \sum_{j=1}^n w_j \right) r_f + \sum_{i=1}^n \left( \mu_i w_j - \frac{1}{2} \sum_{k=1}^n w_j \sigma_{jk} w_k \right) \right].$$

Setting the derivative with respect to  $w_k$  equal to zero, we obtain the equation for the log-optimal portfolio  $\mu_i - r_f - \sum_{j=1}^n \sigma_{ij} w_j = 0$ , which we highlight:

**The log-optimal portfolio** When there is a risk-free asset, the log-optimal portfolio has weights for the risky assets that satisfy

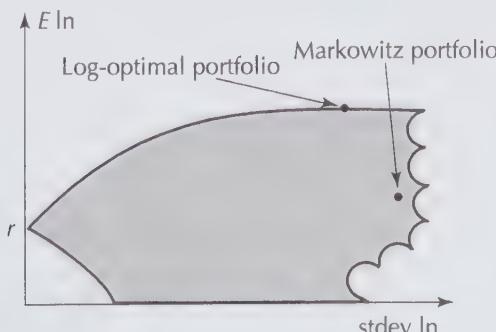
$$\sum_{j=1}^n \sigma_{ij} w_j = \mu_i - r_f \quad (18.6)$$

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

Equation (18.6) is a system of  $n$  linear equations that can be solved for the  $n$  weights.

The efficient frontier with a risk-free asset is shown in Figure 18.5. It should be clear from the figure that most investors will in fact *not* wish to design their strategies to correspond to the log-optimal point. This is because a first-order decrease in standard deviation can be attained with only a second-order sacrifice in expected (log) value by moving slightly leftward along the efficient frontier.

The Markowitz strategy can be defined by using the Markowitz portfolio weights and rebalancing regularly. This strategy will be inefficient with respect to the expected log-variance criterion.



**FIGURE 18.5 The feasible growth rate region.** The Markowitz portfolio is not efficient in the sense of growth.

**Example 18.7 (A single risky asset)** Suppose that there is a single stock with price  $S$  and a riskless bond with price  $B$ . These prices are governed by the equations

$$\begin{aligned}\frac{dS}{S} &= \mu dt + \sigma dz \\ \frac{dB}{B} &= rdt,\end{aligned}$$

where  $z$  is a standard Brownian motion process. The log-optimal strategy will have a weight for the stock given by equation (18.6). In this case that reduces to

$$w = (\mu - r_f)\sigma^2.$$

The corresponding optimal growth rate is

$$\nu_{\text{opt}} = r_f + \frac{(\mu - r_f)^2}{2\sigma^2},$$

and the corresponding variance is

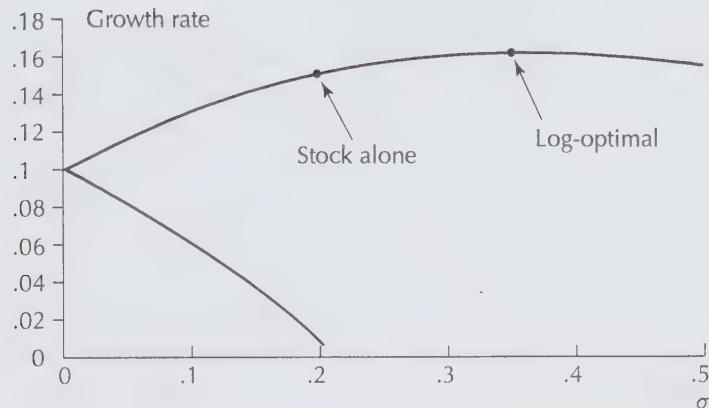
$$\sigma_{\text{opt}}^2 = \frac{(\mu - r_f)^2}{\sigma^2}.$$

Let us consider some numerical values. Suppose that the stock has an expected growth rate of 15% and a standard deviation of 20%. Suppose also that the risk-free rate is 10%. We know that  $\sigma = .20$  and  $\nu = \mu - \frac{1}{2}\sigma^2 = .15$ . This means that  $\mu = .17$ . We find that  $w = 1.75$ , which means that we must borrow the risk-free asset to leverage the stock holding. We also find that the optimal value of  $\nu$  is  $\nu_{\text{opt}} = .10 + (.07)^2/.08 = 16.125\%$ . This is only a slight improvement over the 15% that is obtained by holding the stock alone. Furthermore, the new standard deviation is  $.07/.20 = 35\%$ , which is much worse than that of the stock. The situation is illustrated in Figure 18.6.

The log-optimal strategy does not give much improvement in the expected value, and it worsens the variance significantly. This shows that the log-optimal approach is not too helpful unless there is opportunity to pump between various stocks with high volatility, in which case there can be dramatic improvement.

**Example 18.8 (Three stocks)** Suppose there are three risky stocks with prices governed by the equations

$$\begin{aligned}\frac{dS_1}{S_1} &= .24dt + dz_1 \\ \frac{dS_2}{S_2} &= .20dt + dz_2 \\ \frac{dS_3}{S_3} &= .15dt + dz_3,\end{aligned}$$



**FIGURE 18.6 Feasible region for one stock and a risk-free asset.** The log-optimal strategy gives only modest improvement in growth rate over holding the stock alone, at the expense of a greatly increased standard deviation.

with the covariance of  $d\mathbf{z}$  being

$$\begin{bmatrix} .09 & .02 & .01 \\ .02 & .07 & -.01 \\ .01 & -.01 & .03 \end{bmatrix}$$

The risk-free rate is 10%. We can calculate the corresponding growth rates:  $v_1 = 19.5\%$ ,  $v_2 = 16.5\%$ , and  $v_3 = 13.5\%$ .

Referring to equation (18.6), the log-optimal portfolio weights satisfy the equations

$$\begin{aligned} .09w_1 + .02w_2 + .01w_3 &= .14 \\ .02w_1 + .07w_2 - .01w_3 &= .10 \\ .01w_1 - .01w_2 + .03w_3 &= .05, \end{aligned}$$

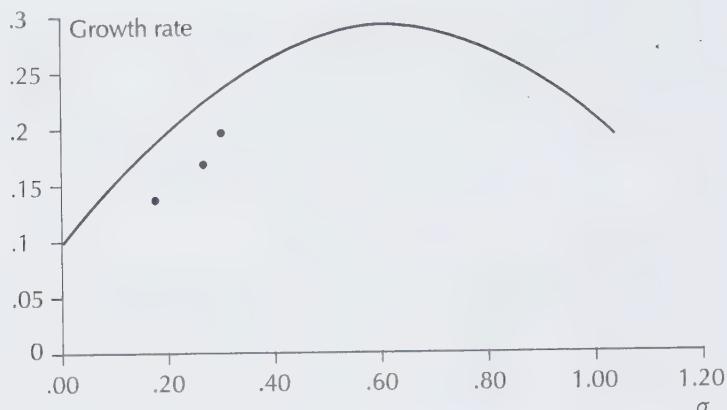
which have solutions

$$w_1 = 1.05$$

$$w_2 = 1.38$$

$$w_3 = 1.78.$$

It follows that  $\mu_{\text{opt}}$  is the corresponding weighted sum of the individual  $\mu$ 's; that is,  $\mu_{\text{opt}} = 1.05 \times .24 + 1.38 \times .20 + 1.78 \times .15 + (1 - 1.05 - 1.38 - 1.78) \times .10 = 0.4742$



**FIGURE 18.7 Boundary points of three-stock example.** The three original stocks together with the risk-free asset define a boundary of points that are optimal with respect to log mean and log variance.

and

$$\begin{aligned}
 \sigma_{\text{opt}}^2 &= \sum_{i,j=1}^3 w_i w_j \sigma_{i,j} \\
 &= .09(1.05)^2 + .02(1.05)(1.38) + .01(1.05)(1.78) + .02(1.38)(1.05) \\
 &\quad + .07(1.38)^2 - .01(1.38)(1.78) + .01(1.78)(1.05) \\
 &\quad - .01(1.05)(1.38) + .03(1.78)^2 \\
 &= 0.3742.
 \end{aligned}$$

Hence  $\sigma_{\text{opt}} = 61.17\%$ . The growth rate is

$$\nu_{\text{opt}} = \mu_{\text{opt}} - \frac{1}{2}\sigma_{\text{opt}}^2 = 28.71\%.$$

Figure 18.7 shows the original three points and a portion of the boundary of the feasible region.

## 18.7 The Log-Optimal Pricing Formula\*

The log-optimal strategy has an important role as a universal pricing asset, and the pricing formula is remarkably easy to derive. As before, we assume that there are  $n$  risky assets with prices each governed by geometric Brownian motion as

$$\frac{dp_i}{p_i} = \mu_i dt + dz_i, \quad i = 1, 2, \dots, n.$$

Since  $E(dz_i) = 0$  for all  $i$ , the covariances  $\sigma_{ij}$  are defined by  $E(dz_i dz_j) = \sigma_{ij} dt$ . There is also a risk-free asset (asset number 0) with rate of return  $r_f$ . Any set of weights  $w_0, w_1, w_2, \dots, w_n$  with  $\sum_{i=0}^n w_i = 1$  defines a portfolio in the usual way. The value of this portfolio will also be governed by geometric Brownian motion. We denote the corresponding covariances of this process with that of asset  $i$  by  $\sigma_{i,\text{port}}$ .

As a special case we denote the log-optimal portfolio by the subscript opt. This portfolio has variance denoted by  $\sigma_{\text{opt}}^2$  and covariance with asset  $i$  denoted by  $\sigma_{i,\text{opt}}$ .

The  $\mu$  of any asset can be recovered from the log-optimal portfolio by evaluating the covariance of the asset with that optimal portfolio. This is essentially a pricing formula because it shows the relation between drift and uncertainty. The pricing formula is stated here (in four different forms):

**Log-optimal pricing formula (LOPF)** *For any stock  $i$  there holds*

$$\mu_i - r_f = \sigma_{i,\text{opt}} \quad (18.7a)$$

$$\nu_i - r_f = \sigma_{i,\text{opt}} - \frac{1}{2}\sigma_i^2. \quad (18.7b)$$

Equivalently, we have

$$\mu_i - r_f = \beta_{i,\text{opt}}(\mu_{\text{opt}} - r_f) \quad (18.8a)$$

$$\nu_i - r_f = \beta_{i,\text{opt}}\sigma_{\text{opt}}^2 - \frac{1}{2}\sigma_i^2, \quad (18.8b)$$

where  $\beta_{i,\text{opt}} = \sigma_{i,\text{opt}}/\sigma_{\text{opt}}^2$ .

**Proof:** The result follows from the equation for the log-optimal strategy (18.6); namely,

$$\mu_i - r_f = \sum_{j=1}^n \sigma_{ij} w_j. \quad (18.9)$$

If  $V$  is the value of the log-optimal portfolio, we have

$$\frac{dV}{V} = \sum_{j=1}^n w_j (\mu_j dt + dz_j).$$

Hence  $\sigma_{i,\text{opt}} = E(dz_i dz_{\text{opt}}) = \sum_{j=1}^n \sigma_{ij} w_j = \mu_i - r_f$ , where the last step is equation (18.9). This gives equation (18.7a). The version of equation (18.7b) follows from  $\nu_i = \mu_i - \frac{1}{2}\sigma_i^2$ .

To obtain the alternative expressions we apply the first pricing formula [equation (18.7a)] to the log-optimal strategy itself, obtaining  $\mu_{\text{opt}} - r_f = \sigma_{\text{opt}}^2$ . Equation (18.8a) follows immediately. The version of equation (18.8b) follows directly from the definition of  $\beta_{i,\text{opt}}$ . ■

According to these formulas the covariance of an asset with the log-optimal portfolio completely determines the instantaneous expected excess return of that asset. Equations (18.7a) and (18.8a), in terms of  $\mu - r_f$ , are easy to remember

because they mimic the CAPM equation. These equations express the excess expected instantaneous return as a single covariance or, in the alternate version, as a beta-type formula.

**Example 18.9 (Three stocks again)** Consider the three stocks of Example 18.8. Let us determine  $\mu_1$  using equation (18.7a). The covariance of  $S_1$  with the log-optimal portfolio is found from

$$E[dz_1(w_1dz_1 + w_2dz_2 + w_3dz_3)] = [1.05 \times .09 + 1.38 \times .02 + 1.78 \times .01]dt = .14dt.$$

Therefore,

$$\mu_1 = r_f + .14 = .24,$$

which is correct since it coincides with the  $\mu_1$  originally assumed.

Equations (18.7b) and (18.8b), in terms of  $v - r_f$ , are perhaps the most relevant equations since  $v$  is the actual observed growth rate. Consider equation (18.8b), which is  $v_i - r_f = \beta_{i,\text{opt}}\sigma_{\text{opt}}^2 - \frac{1}{2}\sigma_i^2$ . For stocks with low volatility (that is, with  $\sigma_i^2$  small), the excess growth rate is approximately proportional to  $\beta_{i,\text{opt}}$ . This parallels the CAPM result. Greater risk leads to greater growth. However, for large volatility the  $-\frac{1}{2}\sigma_i^2$  term comes into play and decreases  $v$ .

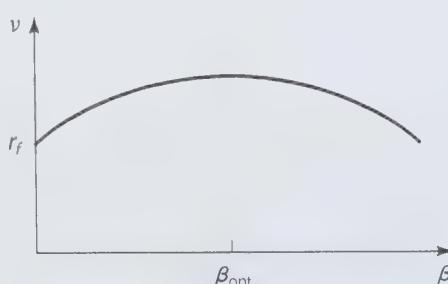
Note in particular that if security  $i$  is uncorrelated with the log-optimal portfolio, its growth rate will be *less* than the risk-free rate. This is because its volatility provides opportunity that a risk-free asset does not.

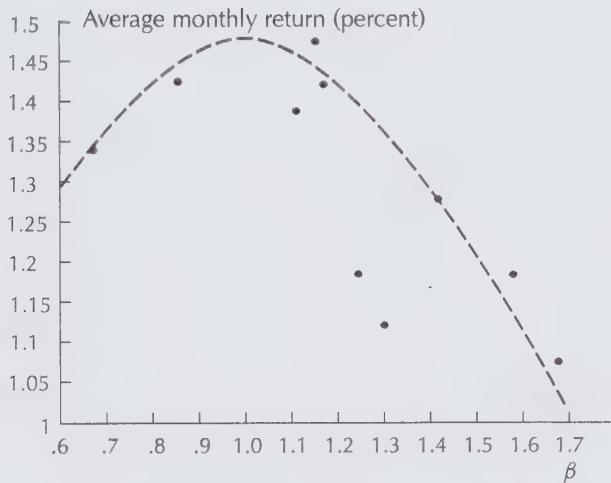
The volatility term implies that the relation between risk and return is quadratic rather than linear as in the CAPM theory. To highlight this quadratic feature, suppose, as may on average be true, that the  $\sigma$  of any stock is proportional to its  $\beta$ ; that is,  $\sigma = \gamma\beta$ , where  $\gamma$  is a constant. Then we find

$$v - r_f = \sigma_{\text{opt}}^2\beta - \frac{\gamma^2\beta^2}{2}.$$

A graph of this function is shown in Figure 18.8. Note that this curve has a different shape than the traditional beta diagram of the CAPM. It is a parabola having a maximum value at  $\beta_{\text{opt}} = \sigma_{\text{opt}}^2/\gamma^2$ .

**FIGURE 18.8 Log return versus beta.**



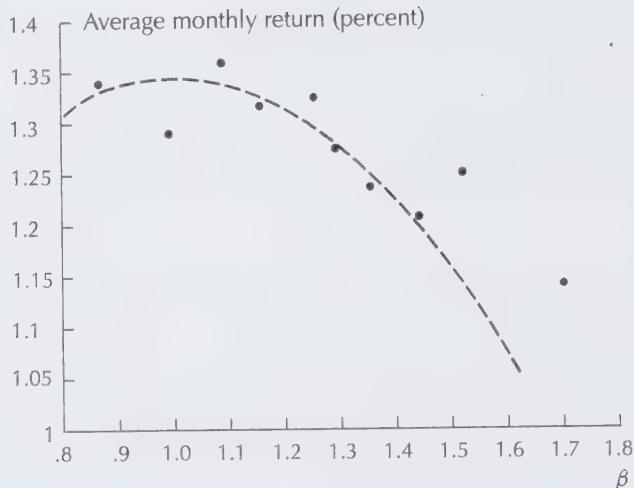


**FIGURE 18.9 Observed return versus  $\beta$  for medium-sized companies.** The data support the conclusion of the LOPF that return is approximately quadratic with respect to  $\beta$  with a peak at around  $\beta = 1$ .

## Market Data

If we were to look at a family of many real stocks, we would not expect them to fall on a single curve like the one shown in Figure 18.8 since the true relationship has two degrees of freedom; namely,  $\beta$  and  $\sigma$ . However, according to the theory discussed, we would expect a scatter diagram of all stocks to fall roughly along such a parabolic curve. We can check this against the results of a famous comprehensive study of market returns which includes many decades of data.<sup>7</sup> The data shown in Figures 18.9 and 18.10 are taken from that study. The figures show annualized return, as computed on a monthly basis, over the period of 1963–1990. Of course the  $\beta$  used in the study is the normal  $\beta$  based on the market return, not on the log-optimal portfolio. This study has been used to argue that the traditional relation predicted by the CAPM does not hold, since the return is clearly not proportional to  $\beta$ . We have drawn a dashed parabola in each figure, which shows that the data *do* support the conclusion that the relation between return and  $\beta$  is roughly quadratic. To put this in perspective, we emphasize that the LOPF is independent of how investors behave. It is a mathematical identity. All that a market study could test, therefore, is whether stock prices really are geometric Brownian motion processes as assumed by the model. Since returns are indeed close to being lognormal, the log-optimal pricing model must closely hold as well.

<sup>7</sup> This is the Fama and French study cited at the end of the chapter. See table I of that reference.



**FIGURE 18.10 Observed return versus  $\beta$  for a cross section of all securities.** The data support the conclusion of the LOPF that return is approximately quadratic with respect to  $\beta$  with a peak at  $\beta = 1$ .

## 18.8 Log-Optimal Pricing and the Black–Scholes Equation\*

The log-optimal pricing formula can be applied to derivative assets, and the resulting formula is precisely the Black–Scholes equation. Hence we obtain a new interpretation of the important Black–Scholes result and see the power of the LOPF. The log-optimal pricing equation is more general than the Black–Scholes equation, since log-optimal pricing applies more generally—not just to derivative assets.

As in the standard Black–Scholes framework, suppose that the price of an underlying asset is governed by the geometric Brownian motion process

$$dS = \mu S dt + \sigma S dz,$$

where  $z$  is a normalized Wiener process. Assume also that there is a constant interest rate  $r$ . Finally, suppose that the price of an asset that is a derivative of the stock is  $y = F(S, t)$  for some (unknown) function  $F$ .

The price  $y$  will fluctuate randomly according to its own Ito process. The equation of this process is given by Ito's lemma as

$$dy(t) = \left( \frac{\partial F}{\partial S} \mu S + \frac{\partial F}{\partial t} + \frac{1}{2} \frac{\partial^2 F}{\partial S^2} \sigma^2 S^2 \right) dt + \frac{\partial F}{\partial S} \sigma S dz. \quad (18.10)$$

If we divide the left side of equation 18.10 by  $y(t)$  and the right side by  $F(S, t)$ , we will have an equation for the instantaneous rate of return of the derivative asset. The first term on the right is then the expected instantaneous rate of return. We can call this  $\mu_{\text{deriv}}$  since it is the  $\mu$  of the derivative asset. Then  $\mu_{\text{deriv}} - r$  must be equal to the covariance of the instantaneous return of the derivative asset with the log-optimal

portfolio. Writing this equation will give the final result. Before we carry this out, let us first find the log-optimal portfolio.

The log-optimal portfolio is a combination of the stock and the risk-free asset. The derivative asset cannot enhance the return achieved by these two assets, since it is by definition a derivative. Therefore the log-optimal portfolio is the combination found in Example 18.7. Specifically, it is the combination in which the weight of the stock is  $w = (\mu - r)/\sigma^2$ .

We can now write the log-optimal pricing formula directly as

$$\frac{1}{F} \left( \frac{\partial F}{\partial S} \mu S + \frac{\partial F}{\partial t} + \frac{1}{2} \frac{\partial^2 F}{\partial S^2} \sigma^2 S^2 \right) - r = \frac{1}{F} \left( \frac{\partial F}{\partial S} \sigma S \right) \left( \frac{\mu - r}{\sigma} \right).$$

The left-hand side is just  $\mu_{\text{deriv}} - r$ , where  $\mu_{\text{deriv}}$  is the expected instantaneous return of the derivative asset. It is found by just copying the first part on the right-hand side of equation (18.10), dividing by  $F$ , and subtracting  $r$ . The right side is the covariance of this derivative asset with the log-optimal portfolio. Since both the derivative and the log-optimal portfolio have prices that are random only through the  $dz$  term, we simply multiply the corresponding coefficients of the instantaneous return equations to evaluate the covariance. The first part is just a copy of the  $dz$  coefficient in equation (18.10) divided by  $F$ , and the second part is the standard deviation of the log-optimal portfolio, as found in Example 18.7.

The equation is simplified by multiplying through by  $F$  and canceling the two identical terms containing  $\mu$ , yielding

$$\frac{\partial F}{\partial t} + \frac{\partial F}{\partial S} r S + \frac{1}{2} \frac{\partial^2 F}{\partial S^2} \sigma^2 S^2 = r F,$$

which is the Black–Scholes equation.

We now have *four* different interpretations of the Black–Scholes equation. The first is a no-arbitrage interpretation, based on the observation that a combination of two risky assets can reproduce a risk-free asset and its rate of return must be identical to the risk-free rate. The second is a backward solution process of the risk-neutral pricing formula. The third is based on the pricing axioms. The fourth is that the Black–Scholes equation is a special case of the log-optimal pricing formula.

## 18.9 Summary

Given the opportunity to invest repeatedly in a series of similar prospects (such as repeated bets on an investment wheel or periodic rebalancing of a stock portfolio), it is wise to compare possible investment strategies relative to their long-term effects on capital. For this purpose, one useful measure is the expected rate of capital growth. If the opportunities have identical probabilistic properties, then this measure is equal to the expected logarithm of a single return. In other words, long-term expected capital growth can be maximized by selecting a strategy that maximizes the expected logarithm of return at each trial; this is the log-optimal strategy.

For bets that pay off either double or nothing, the log-optimal strategy is known as the Kelly rule. It states that you should bet a fraction  $2p - 1$  of your wealth if the probability  $p$  of winning is greater than .5; otherwise, bet nothing.

For stocks, the log-optimal strategy pumps money between volatile stocks by keeping a fixed proportion of capital in each stock, rebalancing each period. This strategy automatically leads, on average, to following the maxim “buy low and sell high.” For stocks, the log-optimal approach is mathematically more tractable in a continuous-time framework than in a discrete-time framework, for in the continuous-time framework explicit formulas can be derived for the log-optimal strategy and the resulting expected growth rate—it is only necessary to solve a quadratic optimization problem. The resulting formula for the expected growth rate clearly shows the source of the pumping effect. Basically: growth rate is  $v = \mu - \frac{1}{2}\sigma^2$ . When assets are combined in proportions, the resulting  $\mu$  is likewise a proportional combination of the individual  $\mu$ 's. However, the resulting  $\sigma^2$  is reduced more than proportionally because it combines individual  $\sigma^2$ 's with squares of the proportionality factors. Therefore the resulting  $v$  is greater than the proportional combination of individual  $v$ 's. Hence  $v$  is pumped up by the reduction in the volatility term.

The growth efficiency proposition states that any long-term investor should evaluate a strategy only in terms of the mean and variance of the logarithm of return. This leads to the concept of an efficient frontier of points on a diagram that shows expected log-return versus standard deviation of log-return. Growth-efficient investors select points on this efficient frontier. This frontier has two extreme points: the log-optimal point and the minimum log-variance point. The two-fund theorem for this framework states that any efficient point is a combination of these two extreme-point portfolios. If there is a risk-free asset, it serves as the minimum log-variance point.

The log-optimal portfolio plays another special role as a pricing portfolio. Specifically, for any asset  $i$ , we find  $\mu_i - r_f = \sigma_{i,\text{opt}}$ . That is, the expected excess instantaneous return of an asset is equal to the covariance of that asset with the log-optimal portfolio. This formula, the log-optimal pricing formula (LOPF), can be transformed to  $v_i - r_f = \beta_{i,\text{opt}}\sigma_{\text{opt}}^2 - \frac{1}{2}\sigma_i^2$ . This shows that the growth rate  $v_i$  tends to increase with  $\beta_{i,\text{opt}}$  as in the CAPM, but it decreases with  $\sigma_i^2$ . Roughly, this leads to security market lines that are quadratic rather than linear. Empirical evidence tends to support this conclusion.

The power of the log-optimal pricing formula (LOPF) is made clear by the fact that the Black–Scholes partial differential equation can be derived directly from the LOPF. However, the LOPF is not limited to the pricing of derivative securities—it is a general result.

## Exercises

- (Simple wheel strategy) Consider a strategy of the form  $(\gamma, 0, 0)$  for the investment wheel. Show that the overall factor multiplying your money after  $n$  steps is likely to be  $(1+2\gamma)^{n/2}(1-\gamma)^{n/2}$ . Find the value of  $\gamma$  that maximizes this factor.
- (How to play the state lottery) In a certain state lottery, people select eight numbers in advance of a random drawing of six numbers. If someone's selections include the six

drawn, they receive a large prize, but this prize is shared with other winners. Victor has discovered that some numbers are “unpopular” in that they are rarely chosen by lottery players. He has computed that by selecting these numbers he has one chance in a million of winning \$10 million for a \$1 lottery ticket. He has odds of 10 to 1 in his favor. Victor’s current wealth is \$100,000, and he wants to maximize the expected logarithm of wealth.

- (a) Should Victor buy a lottery ticket?
  - (b) Victor knows that he can buy a fraction of a ticket by forming a pool with friends. What fraction of a ticket would be optimal?
3. (Easy policy) Show that  $(\frac{1}{2}, \frac{1}{2})$  is the optimal policy for Example 18.2.
4. (A general betting wheel  $\diamond$ ) Consider a wheel with  $n$  sectors. If the wheel pointer lands on sector  $i$ , the payoff obtained is  $r_i$  for every unit bet on that sector. The chance of landing on sector  $i$  is  $p_i$ ,  $i = 1, 2, \dots, n$ . Let  $\alpha_i$  be the fraction of one’s capital bet on sector  $i$ . We require  $\sum_{i=1}^n \alpha_i \leq 1$  and  $\alpha_i \geq 0$  for  $i = 1, 2, \dots, n$ .
  - (a) Show that the optimal growth strategy is obtained by solving
$$\max \sum_{j=1}^n p_j \ln \left( r_j \alpha_j + 1 - \sum_{i=1}^n \alpha_i \right).$$
  - (b) Assuming that  $\alpha_i > 0$  for all  $i = 1, 2, \dots, n$ , show that the optimal values must satisfy
$$\frac{p_k r_k}{r_k \alpha_k + 1 - \sum_{i=1}^n \alpha_i} - \sum_{j=1}^n \frac{p_j}{r_j \alpha_j + 1 - \sum_{i=1}^n \alpha_i} = 0$$

for all  $k = 1, 2, \dots, n$ .

  - (c) Assume that  $\sum_{i=1}^n 1/r_i = 1$ . Show that in this case a solution is  $\alpha_i = p_j$  for  $i = 1, 2, \dots, n$ .
  - (d) For the wheel given in Example 18.5, find the optimal solution and determine the corresponding optimal growth rate.

5. (More on the wheel  $\diamond$ ) Using the notation of Exercise 4, assume that  $\sum_{i=1}^n 1/r_i = 1$ , but try to find a solution where one of the  $\alpha_k$ ’s is zero. In particular, suppose the segments are ordered in such a way that  $p_n r_n < p_i r_i$  for all  $i = 1, 2, \dots, n$ . Then segment  $n$  is the “worst” segment.
 
  - (a) Find a solution with  $\alpha_n = 0$  and all other  $\alpha_i$ ’s positive.
  - (b) Evaluate this solution for the wheel of Example 18.5.

6. (Volatility pumping) Suppose there are  $n$  stocks. Each of them has a price that is governed by geometric Brownian motion. Each has  $v_i = 15\%$  and  $\sigma_i = 40\%$ . However, these stocks are correlated, and for simplicity we assume that  $\sigma_{ij} = .08$  for all  $i \neq j$ . What is the value of  $v$  for a portfolio having equal portions invested in each of the stocks?

7. (The Dow Jones Average puzzle) The Dow Jones Industrial Average is an average of the prices of 30 industrial stocks with equal weights applied to all 30 stocks (but the sum of the weights is greater than 1). Occasionally (about twice per year) one of the 30 stocks splits (usually because its price has reached levels near \$100 per share). When this happens, all weights are adjusted upward by adding an amount  $\epsilon$  to each of them, where  $\epsilon$  is chosen so that the computed Dow Jones Average is continuous.

Gavin Jones’ father, Mr. D. Jones, uses the following investment strategy over a 10-year period. At the beginning of the 10 years, Mr. Jones buys one share of each of the

30 stocks in the Dow Jones average. He puts the stock certificates in a drawer and does no more trading. If dividends arrive, he spends them. If additional certificates arrive due to stock splits, he tosses them in the drawer along with the others. At the end of 10 years he cashes in all certificates. He then compares his overall return, based on the ratio of the final value to the original cost, with the hypothetical return defined as the ratio of the Dow Jones Average now to 10 years ago. He is surprised to see that there is a difference. Which return do you think will be larger? And why? (Ignore transactions costs, and assume that all 30 stocks remain in the average over the 10-year period.) [The difference, when actually measured, is close to 1% per year.]

8. (Power utility) A stock price is governed by

$$\frac{dS}{S} = \mu dt + \sigma dz,$$

where  $z$  is a standardized Wiener process. Interest is constant at rate  $r$ . An investor wishes to construct a constantly rebalanced portfolio to these two assets that maximizes the expected value of his power utility  $U(X) = (1/\gamma)X^\gamma$ ,  $\gamma < 1$ , at all times  $t \geq 0$ . Show that the proportion  $w$  of wealth invested in the risky asset is  $w = (\mu - r)/[(1 - \gamma)\sigma^2]$ . Use the following steps.

- (a) Show that

$$X(t) = X(0)e^{\{rt+w(\mu-r)t-w^2\sigma^2t/2+wn\sigma\sqrt{t}\}},$$

where  $n$  is a normal random variable with mean 0 and variance 1.

- (b) Use  $E(e^{an}) = e^{a^2/2}$  to show that

$$E[U(X(t))] = \frac{1}{\gamma}e^{\gamma\{rt+w(\mu-r)t-w^2\sigma^2t/2+\gamma^2w^2\sigma^2t/2\}}.$$

- (c) Find  $w$ .

9. (Long-term portfolio) You are managing a pension fund with a goal of maximizing the long-term growth rate. There are three assets available. Asset 1 has a risk-free return of 5%. Assets 2 and 3 each are driven by geometric Brownian motion with the following parameter values:  $\mu_2 = 0.1$ ,  $\mu_3 = 0.15$ ,  $\sigma_2 = 0.3$ ,  $\sigma_3 = 0.4$ ,  $\rho = .5$ , where  $\rho$  is the correlation between the two Brownian motions. What are the optimal portfolio weights?

10. (Rebalancing intervals) This exercise explores the sensitivity of log-optimality to the rebalancing frequency. Consider a market consisting of a risk-free asset with zero rate of interest and a stock that over 1 year either increases by a factor  $a > 1$  or decreases by a factor of  $1/a$  with equal probabilities.

- (a) Find the log-optimal portfolio for this situation.  
 (b) Now suppose that the portfolio is rebalanced only every 2 years. Find the log-optimal portfolio in this case.  
 (c) For the case  $a = 2$ , what are the growth rates in parts (a) and (b)? [The growth rate is the yearly expected value of the log of the growth.]

11. (Risk neutral and log optimal) Consider two binomial assets, each with price equal to 1. The first is a stock that at the end of a period pays either 3 or 0 with probabilities  $p$  and  $1 - p$ , respectively. The second asset is risk-free and pays  $R$  at the end of the period.

- (a) A portfolio is defined by investing a portion  $\alpha$  of one's wealth in the stock and a portion  $1 - \alpha$  in the risk-free asset. Find the log-optimal value of  $\alpha$ .

- (b) What is the risk-neutral probability for this situation?  
 (c) What is the value of  $\alpha$  if  $p = q$ ?
- 12.** (Spinning wheel) A game of chance based on a spinning wheel is available that pays  $n$  times money bet in the case of a win and nothing in the case of a loss. A gambler has developed a device by which he may exercise some control over the wheel in such a way that his chance of winning is  $\alpha/n$  instead of  $1/n$ , where  $\alpha > 1$ .
- (a) Starting with \$1.00, how much should the gambler wager to be log optimal?
  - (b) What is the expected value of the log of wealth after a single play?
  - (c) Suppose  $n$  is very large and the gambler uses the strategy  $n$  times. What is the expected value of the log of final wealth as a function of  $n$  as  $n$  goes to infinity?
  - (d) Suppose  $n = 1$  million and you have \$1 million. How much should you bet on the first spin, and how much do you expect earn after a million spins?
- 13.** (Gamma investor) Consider the power utility defined by the function  $F(S) = \frac{1}{\gamma} S^\gamma$ , for  $\gamma \leq 1$ . There are available  $n$  assets, each of which follows geometric Brownian motion, including a risk-free asset with rate of return  $r$ . A portfolio of these leads to a process for the portfolio as  $dP = \mu P dt + \sigma dz$ .
- (a) Show that the gamma investor will seek to maximize  $\mu + \frac{1}{2}(\gamma - 1)\sigma^2$ .
  - (b) Find a pricing formula for  $\mu_i - r$  in terms of  $\text{cov}(P_{\text{opt}}, S_i)$ , where  $P_{\text{opt}}$  is the optimal portfolio of the gamma investor.
- 14.** (Portfolio sensitivity) Suppose there is a stock and a bond governed by the equations
- $$dx = \mu x dt + \sigma x dz$$
- $$dB = rB dt.$$
- It is desired to construct a portfolio of these two securities that gives the maximum expected log of return. However, although  $\sigma$  and  $r$  are known with certainty, the parameter  $\mu$  is uncertain. The best estimate of  $\mu$  is of the form  $\hat{\mu} = \mu + \epsilon$ , where  $\epsilon$  has zero mean and variance  $w^2$ .
- (a) Let  $\alpha$  be the fraction of the wealth invested in the stock at any time and  $1 - \alpha$  be the fraction invested in the bond. If  $\mu$  were known exactly and  $\alpha$  is given, what would be the equation governing the portfolio return  $dP/P$ ?
  - (b) What is the optimal expected log return that would correspond to knowing  $\mu$  exactly? Call it  $R_O$ .
  - (c) If you believe that  $\hat{\mu}$  is the true  $\mu$  and use it to determine  $\alpha$ , what will you think (on average) is the optimal expected log return? Call it  $R_T$ . How does it compare with  $R_O$ ?
  - (d) If you use the estimate  $\hat{\mu}$  to find  $\alpha$ , what is the actual log return that on average you will get? Call it  $R_A$ . How does it compare with  $R_O$ ?
  - (e) Suppose  $\hat{\mu}$  is estimated as equal to the average log return of the stock over  $n$  years. What will be the ratio  $w^2/\sigma^2$ ? If the estimate uses 5 years of data and the expected optimal log return is 15% per year, what are the values of  $R_T$  and  $R_A$ ?
- 15.** (Discrete-time, log-optimal pricing formula) Suppose there are  $n$  assets. Asset  $i$ ,  $i = 1, 2, \dots, n$ , has rate of return  $r_i$  over a single period. There is also a risk-free asset with rate of return  $r_f$ . The log-optimal portfolio over one period has rate of return  $r_0$ , and we define  $P_0 = 1/(1 + r_0)$ .

(a) Derive the pricing formula

$$\bar{r}_i - r_f = -\frac{\text{cov}(r_i, P_0)}{\text{E}(P_0)}.$$

(b) Suppose that over a small period of length  $\Delta t$  the return of asset  $i$  is  $1 + \mu_i \Delta t + n_i \sqrt{\Delta t}$ , where  $n_i$  is a normal random variable with mean 0 and variance  $\sigma_i^2$ . Show that the discrete-time pricing formula in part (a) goes in the limit, as  $\Delta t \rightarrow 0$ , to the continuous-time log-optimal pricing formula given in Section 18.7.

## References

The special advantages of using a logarithmic utility function in situations of repeated investments was initially discovered by Kelly [1] and Latané [2]. The theory was developed more fully by Breiman [3]. See [4] for a good discussion of asymptotic properties. The idea that the expected logarithm and the variance of the logarithm are the only two quantities of importance in long-term behavior was presented in [5]. The fact that the log-optimal portfolio can be used for pricing was presented in [6]. The classic empirical study of security returns is [7].

1. Kelly, J. L., Jr. (1956), "A New Interpretation of Information Rate," *Bell System Technical Journal*, **35**, 917–926.
2. Latané, H. (1959), "Criteria for Choice among Risky Ventures," *Journal of Political Economy*, **67**, 144–155.
3. Breiman, L. (1961), "Optimal Gambling Systems for Favorable Games," Fourth Berkeley Symposium, vol. I, 65–78.
4. Algoet, P. H., and T. M. Cover (1988), "Asymptotic Optimality and Asymptotic Equipartition Properties of Log-Optimum Investment," *Annals of Probability*, **16**, 876–898.
5. Luenberger, D. G. (1993), "A Preference Foundation for Log Mean–Variance Criteria in Portfolio Choice Problems," *Journal of Economic Dynamics and Control*, **17**, 887–906.
6. Long, J. B., Jr. (1990), "The Numeraire Portfolio," *Journal of Financial Economics*, **26**, 29–69.
7. Fama, E. F., and K. R. French (1992), "The Cross-Section of Expected Stock Returns," *Journal of Finance*, **47**, no. 2, 427–465. (See especially Table I.)

# GENERAL INVESTMENT EVALUATION

**C**onceptually, the subject of this chapter is an extension of the fundamental concept of **present value**, the difference being that now the investment is defined by a cash flow process more complex than those studied earlier.

## 19.1 General Present Value

In earlier chapters we learned that present value is defined relative to an existing collection of assets, defining a market. For example, the simplest basic measure of present value in a deterministic single-period system is related to the existing interest rate. Later we found that, in a deterministic setting spanning several periods, present value is dependent on the term structure of interest rates. In a single-period setting with risk, the CAPM provides a price that can be thought of as the present value of a given asset, and this measure depends on the relation of the asset to existing assets in the market.

In this chapter, the market will generally be multiperiod (or in continuous time) and will have assets whose cash flows are random. The cash flows of the assets we wish to value will likewise be multiperiod and random, but they will be more complex than derivatives, in that the randomness of the asset will not necessarily be exactly tied to that of other assets but, instead, may be subject to random factors outside the existing market. In one continuous-time setting, the pricing axioms are used to value non-tradable assets by projection, leading to a generalization of the Black–Scholes equation.

## Projects and Opportunities

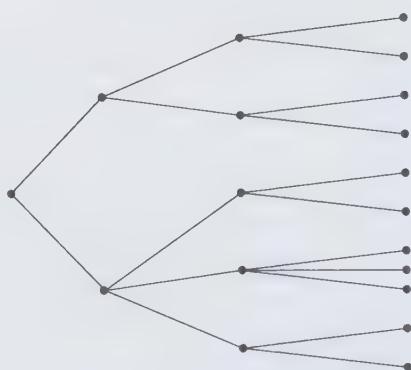
The assets we consider can be regarded as **projects** or **opportunities**, which usually have a finite life. For instance, recall the early examples in Chapter 2 of tree cutting. In one case, tree farming was considered a one-time project, consisting of purchasing land and seedlings and then harvesting the result within a few years. A present-value analysis made perfect sense. On the other hand, we also considered regarding the tree farming business as an ongoing reinvestment in cycles of planting and harvesting, expanding to greater and greater magnitude. This business is more naturally studied as an extension of **internal rate of return** (similar to the investment wheel in Chapter 18). Some projects are standard business ventures of investment followed by a cash flow stream, but they also may be similar to financial derivatives. An example is an option on the revenue of a company rather than on its stock value. This option is not strictly a derivative because there is no underlying security for revenue.

In order to associate a present value to a project or opportunity of the type we have outlined, it is necessary to describe the character of the market environment. In general, this consists of statistical properties of the cash flow process. These can be modeled using the structures we have used before: lattices, trees, or Ito processes. In some cases we will extend these models to include more structure.

## 19.2 Multiperiod Securities\*

We now<sup>1</sup> begin by build a framework for representing securities in a discrete-time multiperiod setting with a finite number of states—a framework that generalizes the discussion of Chapter 11. The basic component of this multiperiod framework is a graph (usually a tree or a lattice) defining a random process of state transitions, as shown in Figure 19.1. The leftmost node represents the initial point of the process

**FIGURE 19.1 State graph.** Each node represents a different state. The graph of this figure is a tree, but in general some nodes may combine.



<sup>1</sup> One may simply scan this section at first and refer to it as necessary.

at time  $t = 0$ . The process can then move to any of its successor nodes at  $t = 1$ . A probability is assigned to each of the arcs. Each probability is greater than or equal to zero, and the sum of the probabilities for arcs emanating from any particular node must be 1.

The nodes of the graph can be thought of as representing various “states of the financial universe.” They might be various weather conditions that would affect agriculture and hence the price of agricultural products. They might be conditions of unemployment that would affect wages and hence profits. Or they might be the various possible prices of gold. The graph must have enough branches to represent the financial problems of interest fully. Particular security processes are defined by assigning values or amounts to the nodes, as discussed next.

## Assets

An asset is defined by a cash flow process, which in turn is defined by assigning a cash flow (or dividend) to each node of the graph. Symbolically, such a cash flow or dividend process is represented by a series of the form  $\delta = (\delta_0, \delta_1, \dots, \delta_T)$ , where each  $\delta_t$  is the cash flow at time  $t$ . The flow  $\delta_t$  is, however, random since it depends on which of the states at time  $t$  actually occurs, so really  $\delta_t$  is a symbol for all possible values at time  $t$ .

Associated with each asset is another process, the price process, which is denoted by  $S = (S_0, S_1, \dots, S_T)$ . The price  $S_t$  represents the price at which the asset would trade after receipt of the cash flow at  $t$ . Again, each  $S_t$  for  $t > 0$  is random since it depends on which node is active at time  $t$ .

The state model can be used to represent several assets simultaneously. Different assets merely correspond to different cash flow and price processes.

The structure of an underlying graph requires some consideration. It is always safest to make this graph a full tree, with no combined nodes. This will assure that any derived quantities can also be accommodated. For computation, on the other hand, we aggressively seek opportunities to combine nodes—perhaps discovering a lattice representation. Then we struggle to keep the nodes from separating, so that we can devise a computationally efficient method of solution.

## Portfolio Strategies

Assume that there are  $n$  assets. Asset  $i$  for  $i = 1, 2, \dots, n$  has the (stochastic) cash flow process  $\delta^i = (\delta_0^i, \delta_1^i, \dots, \delta_T^i)$ . Asset  $i$  also has the stochastic price process  $S_i = (S_0^i, S_1^i, \dots, S_T^i)$ . A **trading strategy** is a portfolio of these assets whose composition may depend on time and on the particular nodes visited. Corresponding to a trading strategy, denoted by  $\theta$ , there is an amount  $\theta_t^i$  of asset  $i$  at time  $t$ , but  $\theta_t^i$  also may depend on the particular state at time  $t$ . In other words, each  $\theta^i = (\theta_0^i, \theta_1^i, \dots, \theta_T^i)$  is itself a process defined on the underlying graph—the process of how much of asset  $i$  is held.

A trading strategy defines a new asset, with an associated cash flow process  $\delta^\theta$ . The cash flows are found from the equation

$$\delta_t^\theta = \sum_{i=1}^n [(\theta_{t-1}^i - \theta_t^i) S_t^i + \theta_{t-1}^i \delta_t^i],$$

where as a convention we put  $\theta_{-1}^i = 0$  for all  $i$ . The first term inside the summation represents the amount of money received at time  $t$  due to changing the portfolio holdings at time  $t$ . The second term is the total dividend received at time  $t$  from the portfolio weights at time  $t-1$ .

As a simple example, consider the trading policy of just buying an asset at time  $t = 0$  for price  $S$  and holding it. This will generate the net cash flow stream  $(-S, \delta_1, \delta_2, \dots, \delta_T)$ .

## Arbitrage

It may be possible to find a strategy that is guaranteed to make money with no cost. Such a strategy is an **arbitrage**. Formally, a trading strategy  $\theta$  is an arbitrage if  $\delta^\theta \geq 0$  and  $\delta^\theta$  is not identically zero. In other words,  $\theta$  is an arbitrage if it generates a dividend process that has at least one positive term and no negative terms. It is easy to imagine an arbitrage, since we have seen many examples in earlier chapters.

## Short-Term Risk-Free Rates

An asset is **short-term risk-free** at time  $t$  if its dividend at time  $t+1$  is  $\delta_{t+1} = 1$  and zero everywhere else. Its price  $S_t$  at time  $t$  gives the discount factor  $d_t = S_t$ . Purchase of this security at time  $t$  yields the cash flow process  $(0, 0, \dots, -S_t, 1, 0, \dots, 0)$ . If there is no such underlying asset, it may be possible to construct one synthetically with a trading strategy. In either case, we say that short-term risk-free borrowing exists. We define the risk-free return as  $R_{t,t+1} = 1/d_t$ .

## 19.3 Risk-Neutral Pricing

We now turn to one of the main themes emphasized throughout the book: risk-neutral pricing. We shall use a discrete-time approach. We assume throughout this section that short-term risk-free rates exist for all periods.

Assume that there are  $n$  assets defined on the underlying state process graph. From these assets, new assets can be constructed by using trading strategies. We say that risk-neutral probabilities exist if a set of risk-neutral probabilities can be assigned to the arcs of the graph such that the price of any asset or any trading policy satisfies

$$S_t = \frac{1}{R_{t,t+1}} \hat{E}_t(S_{t+1} + \delta_{t+1}) \quad (19.1)$$

for every  $t = 0, 1, 2, \dots, T - 1$  and where  $\hat{E}_t$  denotes expectation at time  $t$  with respect to the risk-neutral probabilities.

This definition applies only one period at a time, and it is expressed in a backward fashion. It gives  $S_t$  as a function of the reachable values of  $S_{t+1}$  and  $\delta_{t+1}$ .

We cannot assume that risk-neutral probabilities exist for the particular set of assets in our collection. However, as one might suspect, we can guarantee the existence of risk-neutral probabilities when the prices of the original assets are consistent in a way that makes arbitrage impossible. This is the content of the following theorem, which follows immediately from our earlier result in Chapter 11 on risk-neutral pricing because the risk-neutral pricing formula (19.1) is a single-period formula.

**Existence of risk-neutral probabilities** *Suppose a set of  $n$  assets is defined on a state process. Suppose that from these assets, short-term risk-free borrowing is possible at every time  $t$ . Then there are risk-neutral probabilities such that the prices of trading strategies with respect to these assets are given by the risk-neutral pricing formula*

$$S_t = \frac{1}{R_{t,t+1}} \hat{E}_t(S_{t+1} + \delta_{t+1})$$

*if and only if no arbitrage is possible.*

**Proof:** We already have all the elements. It is clear that risk-neutral pricing implies that no arbitrage is possible. This was shown in Section 16.3 for a short rate lattice, and the proof carries over almost exactly.

It remains to be shown that if no arbitrage is possible, then there are risk-neutral probabilities. However, if no arbitrage is possible over the  $T$  periods, certainly no arbitrage is possible over the single period at  $t$ , starting at a given node. It was shown in Chapter 11 that this implies that risk-neutral probabilities exist for the arcs emanating from that node. Since this is true for all nodes at all times  $t$ , we obtain a full set of risk-neutral probabilities. ■

The risk-neutral pricing formula (19.8) can be written in a nonrecursive form as

$$S_t = \hat{E}_t \left( \sum_{s=t+1}^T \frac{\delta_s}{R_{ts}} \right), \quad (19.2)$$

where now  $\hat{E}_t$  denotes expectation of all future quantities starting at the known state at time  $t$ . This formula expresses  $S_t$  as a discounted risk-neutral evaluation of the entire remaining cash flow stream. It has the nice interpretation of generalizing the familiar present value formula used for deterministic cash flow streams. However, this form is not convenient for calculation because the quantity  $R_{ts}$  generally requires a full tree representation. (See Exercise 2.) There are cases where the result simplifies, of course, such as when interest rates are deterministic.

## 19.4 Optimal Pricing

According to the foregoing theorem, risk-neutral probabilities exist if there is no opportunity for arbitrage among the available assets. The theorem does not say that these probabilities are unique, and, in general, they are not.

If the assets span the degrees of freedom in the underlying graph, as is the case of two assets on a binomial lattice, then the risk-neutral prices *are* unique. If they do not span, as in the case of two assets on a trinomial lattice, there will be additional degrees of freedom, and the risk-neutral probabilities are not unique.

When there are extra degrees of freedom, a specific set of risk-neutral probabilities can be defined by introducing a utility function  $U$ , measuring the utility of the final wealth level, and finding the trading policy that maximizes the expected value of  $U(X_T)$ . This optimal trading policy will imply a set of risk-neutral prices in a manner similar to that for the single-period case discussed in Chapter 11.

We shall limit our consideration to utility functions that have a separation property (as was done in Section 18.4). The separation property holds for the logarithm, and it also holds for the power utility function  $U(X_T) = (1/\gamma)X_T^\gamma$ . When the separation property holds, the multiperiod case reduces to a series of single-period problems, all having the same form of utility function. This greatly simplifies the necessary calculations (although most of the general conclusions hold for other utility functions).

### The Single-Period Problem

Recall that there are  $n$  assets. The single-period problem at time  $t$ , and at a specific node at that time, is to select amounts  $\theta_t^i$  for  $i = 1, 2, \dots, n$  of the  $n$  assets, forming a portfolio. We wish to maximize the expected utility of the value of this portfolio at  $t+1$  subject to the condition that the total cost of the portfolio at time  $t$  is 1. Hence we seek  $\theta_t^i$ 's to solve

$$\underset{\theta_t}{\text{maximize}} \quad E_t[U(X_{t+1})] \quad (19.3)$$

$$\text{subject to} \quad \sum_{i=1}^n \theta_t^i S_t^i = 1 \quad (19.4)$$

$$\sum_{i=1}^n \theta_t^i (S_{t+1}^i + \delta_{t+1}^i) = X_{t+1}. \quad (19.5)$$

The expectation is taken with respect to the actual probabilities of successor nodes. If there are  $K$  such nodes, we denote these probabilities by  $p_1, p_2, \dots, p_K$ . Given amounts  $\theta_t^i$ ,  $i = 1, 2, \dots, n$ , the value of next-period wealth  $X_{t+1}$  depends on the particular successor node  $k$  that occurs. The objective function can be written as  $\sum_{k=1}^K p_k U(X_{t+1})_k$ , where  $U(X_{t+1})_k$  denotes the value of  $U(X_{t+1})$  at node  $k$ .

Using the construction of equation (11.16), a set of risk-neutral probabilities can be found from the solution. Specifically, the risk-neutral probabilities are

$$q_k = \frac{p_k U'(X_{t+1}^*)_k}{\sum_{k=1}^K p_k U'(X_{t+1}^*)_k}, \quad (19.6)$$

where  $X_{t+1}^*$  is the optimal (random) value of next-period wealth. If the utility function  $U$  is increasing,  $U'(X_{t+1}^*)_k$  will be positive, and hence all the  $q_k$ 's will be positive. (It is assumed that  $X_{t+1}^* > 0$ .)

These risk-neutral probabilities can be used to price any asset using the general formula

$$S_t = \frac{\hat{E}_t(S_{t+1} + \delta_{t+1})}{R_{t,t+1}},$$

which takes the specific form

$$S_t = \frac{\sum_{k=1}^K q_k (S_{t+1} + \delta_{t+1})_k}{R_{t,t+1}}.$$

## Applications

If this method is used to find a set of risk-neutral probabilities when there are more states than basic assets, the risk-neutral probabilities will depend on the choice of utility function and wealth. The variations in the risk-neutral probabilities will not affect the prices of the original assets, but will lead to variations in the prices assigned to other (new) assets. The price assigned to a new asset this way is such that an individual with the given utility function will not choose to include that asset in the optimal portfolio (either long or short). That is, it is a zero-level price for that investor.

**Example 19.1 (A 5-month call)** As a specific example let us consider the 5-month call option studied in Example 14.3. The underlying stock had  $S(0) = \$62$ ,  $\mu = .12$ , and  $\sigma = .20$ . The risk-free rate is  $r = 10\%$  per annum, and the strike price of the option is  $K = \$60$ .

For better accuracy we use a trinomial lattice with 1-month periods. To match the parameters of the stock, we decide on the trinomial parameters  $u = 1.1$ ,  $d = 1/u$ , and the middle branch has a multiplicative factor of 1. To find the real probabilities we must solve the equations that correspond to: (1) having the probabilities sum to 1, (2) matching the mean, and (3) matching the variance. These equations, first given in Section 15.7, are

$$\begin{aligned} p_1 + p_2 + p_3 &= 1 \\ up_1 + p_2 + dp_3 &= 1 + \mu \Delta t \\ u^2 p_1 + p_2 + d^2 p_3 &= \sigma^2 \Delta t + (1 + \mu \Delta t)^2. \end{aligned}$$

They have solution  $p_1 = .228$ ,  $p_2 = .632$ , and  $p_3 = .140$ .

Now that the lattice parameters are fixed, we must solve one step of the optimal-portfolio problem in order to find suitable risk-neutral probabilities. For this purpose we chose a logarithmic utility function. Hence we solve the problem

$$\max_{\alpha} E\{\ln[\alpha R + (1 - \alpha)R_0]\},$$

where  $R$  is the random return of the stock over one period and  $R_0$  is the risk-free return. Written out in detail this is

$$\max\{p_1 \ln[\alpha u + (1 - \alpha)R_0] + p_2 \ln[\alpha + (1 - \alpha)R_0] + p_3 \ln[\alpha d + (1 - \alpha)R_0]\}.$$

This has optimal solution  $\alpha = .505$ . The corresponding risk-neutral probabilities are then readily found from (19.6) to be

$$q_1 = \frac{p_1}{\alpha u + (1 - \alpha)R_0} c \quad (19.7)$$

$$q_2 = \frac{p_2}{\alpha + (1 - \alpha)R_0} c \quad (19.8)$$

$$q_3 = \frac{p_3}{\alpha d + (1 - \alpha)R_0} c, \quad (19.9)$$

where  $c$  is the normalizing constant. When normalized the values are  $q_1 = .218$ ,  $q_2 = .635$ , and  $q_3 = .148$ .

**FIGURE 19.2 Log-optimal pricing of a 5-month call option using a trinomial lattice.** The upper lattice contains the possible stock prices. The lower lattice is found by risk-neutral valuation using inferred probabilities.

					99.85
					90.77
				82.52	82.52
			75.02	75.02	75.02
		68.20	68.20	68.20	68.20
	62.00	62.00	62.00	62.00	62.00
		56.36	56.36	56.36	56.36
			51.24	51.24	51.24
				46.58	46.58
					42.35
					38.50
					39.85
					31.27
				23.51	23.02
			16.51	16.01	15.52
		10.43	9.85	9.27	8.70
	5.8059	5.20	4.56	3.85	3.03
		1.92	1.43	.92	.43
			.26	.09	.00
				.00	.00
				.00	.00
					.00

					30.77
					22.52
				15.52	15.02
			8.70		8.20
					2.00
					.00
					.00
					.00

With these values in hand it is possible to proceed through the lattice in the normal backward solution method. The results are shown in Figure 19.2. The price obtained is \$5.8059, which is very close to the Black–Scholes value of \$5.80.

## 19.5 The Double Lattice

The starting point for general investment analysis as presented in this chapter is a graph that represents a family of asset processes. How can we construct such a graph to embody the characteristics of each asset and the relations between assets? Clearly, this construction may be quite complex.

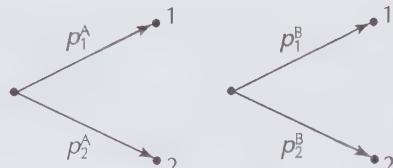
This section shows how a graph for two risky assets can be constructed by combining the separate representations for each asset. Specifically, two binomial lattices are combined to produce a double lattice that faithfully represents both assets.

Suppose that we have two assets A and B, each represented by a binomial lattice. Each has up and down factors and probabilities, but movements in the two may be correlated. A representation of one time period is shown in Figure 19.3.

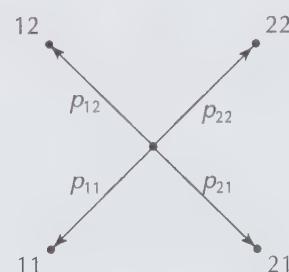
The combination of these two lattices is really a lattice with *four* branches at each time step. It is most convenient to use double indexing for this new combined lattice; call the nodes 11, 12, 21, and 22. The first index refers to the first lattice and the second to the second. We define the transition probabilities as  $p_{11}$ ,  $p_{12}$ ,  $p_{21}$ , and  $p_{22}$ , respectively. A picture of the combined lattice is shown in Figure 19.4. Here the center node is the node at an initial time, and the four outer nodes are the four possible successors.

Suppose the lattice for stock A has node factors  $u^A$  and  $d^A$  with probabilities  $p_1^A$  and  $p_2^A$ , respectively; and the lattice for stock B has node factors  $u^B$  and  $d^B$  with probabilities  $p_1^B$  and  $p_2^B$ . If the covariance of the logarithm of the two return factors

**FIGURE 19.3 One step of two separate lattices.** Their movements may be correlated.



**FIGURE 19.4 Nodes of the combination.** There are four possible successor nodes from the central node.



$\sigma_{AB}$  is known, we may select the probabilities of the double lattice to satisfy

$$p_{11} + p_{12} = p_1^A$$

$$p_{21} + p_{22} = p_2^A$$

$$p_{11} + p_{21} = p_1^B$$

$$(p_{11} - p_1^A p_1^B) U^A U^B + (p_{12} - p_1^A p_2^B) U^A D^B$$

$$+ (p_{21} - p_2^A p_1^B) D^A U^B + (p_{22} - p_2^A p_2^B) D^A D^B = \sigma_{AB},$$

where  $U^A = \ln u^A$ ,  $D^A = \ln d^A$ ,  $U^B = \ln u^B$ , and  $D^B = \ln d^B$ .

A special case is when the covariance is zero, corresponding to independence of the two asset returns. In that case it follows that the appropriate lattice probabilities are  $p_{11} = p_1^A p_1^B$ ,  $p_{12} = p_1^A p_2^B$ ,  $p_{21} = p_2^A p_1^B$ , and  $p_{22} = p_2^A p_2^B$ .

Linearity of the set of equations for the  $p_{ij}$ 's, implies that if  $\sigma_{AB} \neq 0$ , each of the unknown probabilities will change from the independent case by the addition or subtraction of a multiple of  $\sigma_{AB}$ . Based on that observation it is easily shown that

$$p_{11} = p_1^A p_1^B + \bar{c}$$

$$p_{12} = p_1^A p_2^B - \bar{c}$$

$$p_{21} = p_2^A p_1^B - \bar{c}$$

$$p_{22} = p_2^A p_2^B + \bar{c},$$

where

$$\bar{c} = \frac{\text{COV}_{AB}}{(U^A - D^A)(U^B - D^B)}.$$

**Example 19.2 (Two nice stocks)** Consider two stocks with identical binomial lattice representations of  $u = 1.3$ ,  $d = .9$ , and  $p_u = .6$ ,  $p_d = .4$ . Assume also that they have a correlation coefficient of  $\rho = .3$ . Let us find the double lattice representation.

Let  $S_A$  and  $S_B$  be the random values of the two stocks after one period when initiated at unity. We have

$$E(\ln S_A) = E(\ln S_B) = .6 \times \ln 1.3 + .4 \times \ln .9 = .11527$$

$$\sigma^2 = \text{var}(\ln S_A) = \text{var}(\ln S_B) = .6(\ln 1.3)^2 + .4(\ln .9)^2 - .11527^2 = .03245$$

$$\text{cov}_{AB} = .3\sigma^2 = .009736.$$

We find  $\bar{c} = \text{cov}_{AB}/(\ln 1.3 - \ln .9)^2 = .072$ . Hence we immediately obtain the solution

$$p_{11} = .6 \times .6 + .072 = .432$$

$$p_{12} = .6 \times .4 - .072 = .168$$

$$p_{21} = .4 \times .6 - .072 = .168$$

$$p_{22} = .4 \times .4 + .072 = .232.$$

## 19.6 Pricing in a Double Lattice

The double lattice construction does provide a valid representation of the two assets, but there is a problem. When a risk-free asset is adjoined, we have four nodes, but only *three* assets: the two risky assets and the risk-free asset. There is an extra degree of freedom. Therefore the risk-neutral probabilities are not completely specified as they are in the two original small lattices. We must find a way to pin down that extra degree of freedom in the definition of the risk-neutral probabilities.

One way to specify risk-neutral probabilities is to introduce a utility function, as in the previous section. Different utility functions may lead to different risk-neutral probabilities, but it turns out that under certain conditions a fourth relation holds independently of the particular utility function.

Let us introduce a utility function  $U$ . We determine the risk-neutral probabilities by maximizing expected utility. Denote the optimal value of wealth at the next time point, at node  $ij$ , by  $X_{ij}^*$ ; and, correspondingly, define  $U'_{ij} = U'(X_{ij}^*)$ . Then the risk-neutral probabilities are, from (19.13),

$$q_{ij} = \frac{p_{ij} U'_{ij}}{\sum_{k,l=1}^2 p_{kl} U'_{kl}} \quad (19.10)$$

for  $i, j = 1, 2$ . If the utility function  $U$  is strictly increasing, then the risk-neutral probabilities are strictly positive.

In certain special cases there will be a relation among the  $q_{ij}$ 's that will supply the additional relation needed to make them unique. Two of those cases are spelled out in the following theorem:

**Ratio theorem** Suppose the  $q_{ij}$ 's are determined by (19.10). Then the relation

$$\frac{q_{11}q_{22}}{q_{12}q_{21}} = \frac{p_{11}p_{22}}{p_{12}p_{21}}$$

holds if either of the following two conditions is satisfied:

- (a) One of the original assets appears at zero level in the optimal portfolio.
- (b) The time  $\Delta t$  between periods is vanishingly small.

**Proof:** We shall prove that under either condition  $U'_{11}U'_{22} = U'_{12}U'_{21}$ . This fact will then lead to the final conclusion.

- (a) Suppose that asset A has zero level in the optimal portfolio. Then changes in asset A do not influence  $U'$ . Hence  $U'_{11} = U'_{21}$  and  $U'_{12} = U'_{22}$ . Therefore

$U'_{11} U'_{22} = U'_{12} U'_{21}$ . Clearly, the same result holds if asset B has zero level in the optimal portfolio.

- (b) Now, as a second case, assume that  $\Delta t$  is small. At the optimal portfolio we may write  $X_{ij} = (\tilde{R}_i^A + \tilde{R}_j^B + \tilde{R}^0)X_t$ , where the terms  $\tilde{R}_i^A, \tilde{R}_j^B$ , and  $\tilde{R}^0$  are the returns in the portfolio that correspond to the risky asset A, the risky asset B, and the risk-free asset, respectively. For small  $\Delta t$  the return over one period must be close to 1. Hence,

$$\begin{aligned}\tilde{R}_i^A + \tilde{R}_j^B + \tilde{R}^0 &= 1 + r_i^A + r_j^B + r^0 \\ &\approx (1 + r_i^A)(1 + r_j^B)(1 + r^0),\end{aligned}$$

where  $r_i^A, r_j^B$ , and  $r^0$  are small. This approximation carries over to  $U'$  as well, giving

$$\begin{aligned}U'_{ij} &= U'[(1 + r_i^A + r_j^B + r^0)X_t] \\ &\approx U'(X_t) + U''(X_t)(r_i^A + r_j^B + r^0)X_t \\ &\approx U'(X_t)(1 + \gamma r_i^A)(1 + \gamma r_j^B)(1 + \gamma r^0),\end{aligned}$$

where

$$\gamma = \frac{U''(X_t)X_t}{U'(X_t)}.$$

This product form for  $U'_{ij}$  implies that

$$U'_{11} U'_{22} = U'_{12} U'_{21}.$$

Under condition (a) or (b) we have  $U'_{11} U'_{22} = U'_{12} U'_{21}$ . We then compute

$$\frac{q_{11}q_{22}}{q_{12}q_{21}} = \frac{p_{11}U'_{11}p_{22}U'_{22}}{p_{12}U'_{12}p_{21}U'_{21}} = \frac{p_{11}p_{22}}{p_{12}p_{21}}. \blacksquare$$

An important special case of the two lattice construction is where the two original lattices are independent. In that case  $p_{11} = p_1^A p_1^B$ ,  $p_{12} = p_1^A p_2^B$ ,  $p_{21} = p_2^A p_1^B$ , and  $p_{22} = p_2^A p_2^B$ . It follows by direct substitution that

$$\frac{p_{11}p_{22}}{p_{12}p_{21}} = 1.$$

Then if either of the conditions of the ratio theorem is satisfied,

$$\frac{q_{11}q_{22}}{q_{12}q_{21}} = 1$$

and from this it can be shown that the original two lattices are independent with respect to the risk-neutral probabilities as well as with respect to the original probabilities. That is, independence with respect to original probabilities implies independence with respect to risk-neutral probabilities.<sup>2</sup>

Now let us return to our original problem. In the double lattice we have four successor nodes but only three assets. For small  $\Delta t$ , the ratio formula gives the fourth relation required to determine a set of four risk-neutral probabilities.

An important special case of the two-lattice situation is that where one of the lattices is a short rate lattice for interest rates. This case can be treated by the same technique, as illustrated by the Simplico gold mine example that follows.

**Example 19.3 (Double stochastic Simplico gold mine)** Consider a 10-year lease on the Simplico mine. In evaluating this lease we recognize that the price of gold and the interest rate are *both* stochastic, but we will assume that they are independent.

Recall that up to 10,000 ounces of gold can be extracted from this mine each year at a cost of \$200 per ounce. The price of gold is initially \$400 per ounce and fluctuates according to a binomial lattice that has an up factor of  $u = 1.2$  and a down factor of  $d = .9$ . The price obtained for sale of the gold produced in a year is assumed to be the gold price at the beginning of the year, but the cash flow occurs at the end of the year.

In this version of the problem we assume that the term structure of interest rates is governed by a short rate lattice. The initial short rate is 4%, and the lattice is a simple up-down model with  $u' = 1.1$  and  $d' = .9$ . The risk-neutral probabilities are given as .5. We shall use the small  $\Delta t$  approximation to assert that the result of the ratio theorem applies. Then since the gold price fluctuations and the short rate fluctuations are independent of each other, we conclude that the risk-neutral probabilities are also independent. Hence the actual probabilities are irrelevant for pricing purposes.

We solve this problem by constructing a double lattice. Each node of this lattice represents a combination  $(g, r)$  of gold price  $g$  and short rate  $r$ . Each of these nodes is connected to four neighbor nodes with values  $(ug, u'r)$ ,  $(ug, d'r)$ ,  $(dg, u'r)$ , and  $(dg, d'r)$ . The risk-neutral probabilities of these arcs are just the product of the risk-neutral probabilities for movement in the two elementary lattices. For interest rates, these are each .5. For gold, the probability of an up move is  $q_u = (1 + r - d)/(u - d)$ , where  $r$  is the (current) short rate. Hence the four probabilities for the double lattice, corresponding to arcs leading to the nodes listed, are  $q_{11} = .5q_u$ ,  $q_{12} = .5q_u$ ,  $q_{21} = .5(1 - q_u)$ , and  $q_{22} = .5(1 - q_u)$ .

The double lattice can be set up as a series of 10 two-dimensional arrays. Each array contains the possible  $(g, r)$  pairs for that period. The arrays are then linked by the risk-neutral pricing formula. This formula simply multiplies the values at each of the four successor nodes by their risk-neutral probabilities, adds those plus the cash flow for the end of the year, and discounts the sum using the current short rate. The

<sup>2</sup> Briefly: Let  $\mathbf{Q}$  be the  $2 \times 2$  matrix with components  $[q_{ij}]$ . Then the invariance condition says that  $\mathbf{Q}$  is singular, which means  $\mathbf{Q} = \mathbf{ab}^T$  for some  $2 \times 1$  vectors  $\mathbf{a}$ ,  $\mathbf{b}$ . Normalization makes both of these vectors have components that sum to 1; and these define the individual probabilities.

Period 9											<i>r</i>
<i>g</i>	0.0155	0.0189	0.0231	0.0283	0.0346	0.0423	0.0517	0.0631	0.0772	0.0943	<i>r</i>
2063.91	18.355	18.293	18.217	18.126	18.016	17.883	17.724	17.532	17.304	17.033	
1547.93	13.274	13.229	13.174	13.108	13.029	12.933	12.817	12.679	12.514	12.318	
1160.95	9.463	9.431	9.392	9.345	9.288	9.220	9.137	9.039	8.921	8.781	
870.71	66.048	6.582	6.555	6.523	6.483	6.435	6.378	6.309	6.227	6.129	
653.03	4.461	4.446	4.428	4.406	4.379	4.347	4.308	4.261	4.206	4.140	
489.78	28.535	2.844	2.832	2.818	2.801	2.780	2.755	2.726	2.690	2.648	
367.33	1.648	1.642	<b>1.635</b>	<b>1.627</b>	1.617	1.605	1.591	1.574	1.553	1.529	
275.50	7.435	0.741	<b>0.738</b>	<b>0.734</b>	0.730	0.724	0.718	0.710	0.701	0.690	
206.62	0.065	0.065	0.065	0.064	0.064	0.064	0.063	0.062	0.061	0.061	
154.97	0	0	0	0	0	0	0	0	0	0	

Period 8											<i>r</i>
<i>g</i>	0.0172	0.0210	0.0257	0.0314	0.0384	0.0470	0.0574	0.0702	0.0857		<i>r</i>
1719.93	29.917	29.812	29.685	29.531	29.345	29.121	28.852	28.529	28.144		
1289.95	21.463	21.390	21.301	21.194	21.064	20.907	20.719	20.493	20.224		
967.46	32.925	15.073	15.013	14.941	14.853	14.747	14.619	14.466	14.283		
725.59	21.784	10.336	10.297	10.251	10.194	10.126	10.044	9.946	9.828		
544.20	14.492	6.782	6.760	6.733	6.701	6.661	6.613	6.555	6.486		
408.15	9.058	4.118	4.107	4.095	4.080	4.062	4.040	4.013	3.980		
306.11	4.123	2.119	<b>2.118</b>	2.117	2.115	2.113	2.110	2.106	2.100		
229.58	1.900	0.620	0.626	0.633	0.641	0.651	0.662	0.675	0.690		
172.19	0.025	0.026	0.026	0.027	0.028	0.030	0.031	0.033	0.035		

**FIGURE 19.5 Arrays for two periods of the Simplico gold mine.** Each mode at period  $k$  has four successor nodes at period  $k+1$ , as indicated by the corresponding bold numbers. Values are in millions of dollars.

values at time 10 are all zero. Figure 19.5 shows the values at the nodes for time periods 9 and 8.

The first column shows the possible  $g$  values and the first row shows the possible  $r$  values. The entries in the main arrays are the corresponding values (in millions of dollars) of the lease. A node in the period 8 array is found from four nodes in the period 9 array, as illustrated in the figure.

Working backward this way we find an array with just one node at period zero, having a value of \$22.2551 million dollars.

## 19.7 Investments with Private Uncertainty

Suppose a project requires an initial cash outlay and will produce an uncertain cash flow at the end of one year. Suppose also that the uncertainty consists of both **private** uncertainty and **market** uncertainty. Basically, market uncertainty can be replicated with market participation, whereas private uncertainty cannot. For example, the cash

flow of a gold mine lease depends both on the market uncertainty of gold prices and on the private uncertainty of how much gold is in the yet unexplored veins.

One way to assign a value to such a project is to make believe that the project value is a price, and then set the price so that you would be indifferent between either purchasing a small portion of the project or not. This is termed **zero-level pricing** since you will purchase the project at zero level. Of course, it is assumed that you have the option to purchase other assets, including at least a risk-free security with total return  $R$ .

If there is only private uncertainty the zero-level price is just the discounted expected value of the project (using actual probabilities). It cannot be priced any lower, for then you would want to purchase a small amount of it. Likewise, it cannot be priced any higher, or you would want to sell (short) some of it. The value is therefore

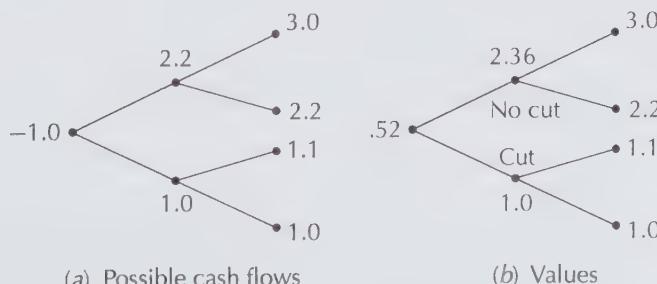
$$V = c_0 + \frac{1}{R} E(c_1)$$

where  $c_0$  and  $c_1$  are the initial and final cash flows, respectively.

Notice that this is somewhat different than the formula for the price of market assets. Market assets already have prices, and you will likely want to include them in your portfolio at a nonzero level (either long or short).

**Example 19.4 (When to cut a tree)** Suppose that we can grow trees (for lumber). The trees grow randomly, and the cash flows associated with harvest after 1 year or after 2 years are shown by the (diagram) tree on the left side of Figure 19.6. Each arc of the tree has a probability of .5. The uncertainty is private because the tree growth depends only on local weather conditions and is not related to market variables.

The initial cash flow of  $-1.0$  must be paid to carry out the project. The cash flow figures shown at the end of the period are those that will be received if the trees are cut after 1 year. Likewise, the final values shown are the cash flows that will be received if the trees are not cut until after 2 years. We wish to evaluate this project, assuming that the interest rate is constant at 10%. To do so we will need to determine the best strategy for cutting the trees.



**FIGURE 19.6 When to cut a tree.** (a) Cash flow generated at a node if the trees are cut at that point. (b) Value at a node and best policy. (a) Possible cash flows (b) Values

We use the zero-level pricing method, and since there is more than one period, we work backward in the usual fashion. The expected value of the top two nodes at the last time period is 2.6. Discounted by 10% this is a value of 2.36. Since this is higher than 2.2, this is the best value that can be attained if we arrive at the upper node after 1 year. We record this optimal value on the values diagram in Figure 19.6(b). We also place a notation near that node that we should not cut the trees if we arrive there. Likewise, the expected value of the bottom two nodes at the last time period is 1.05. Discounted, this is .95, which is less than 1.0, so we would assign 1.0 at the next backward node in the values diagram, and place a notation there that we should cut the trees if we arrive at that node. The expected value of these two optimal one-period values is  $.5(2.36 + 1.0) = 1.68$ , which discounted is 1.52. Hence the overall value is .52.

## General Approach

The preceding result concerning zero-level pricing of projects with private uncertainty can be generalized to projects that are characterized as having both private uncertainty and market uncertainty. The private uncertainties include such things as unknown production efficiency (due to new production processes), uncertainty in resources (such as the amount of oil in an oil field), uncertainty of outcome (as in a research and development project), and a component of the price uncertainty of commodities for which there is no liquid market (such as the future price of an isolated piece of land). Market uncertainties are those associated with prices of traded commodities and other assets.

Formally, suppose that the states of the world are factored into two parts: a market component and a nonmarket (private) component. A general state (or node in the state graph) therefore can be written as  $(s_t^m, s_t^n)$  corresponding to the market and nonmarket components at time  $t$ . For simplicity we assume that these two components are statistically independent.

From a given state there are various successor states. In the lattice framework we index the successor market states (which are nodes in the lattice) by  $i$  and the nonmarket nodes by  $j$ . The probability of the  $i$ th market node is  $p_i^m$  and the probability of the  $j$ th nonmarket node is  $p_j^n$ . Since the two components are independent, the probability of  $i$  and  $j$  together is  $p_{ij} = p_i^m p_j^n$ . We are now in the situation of a double tree or double lattice.

We also assume that the market portion of the system is complete in the sense that there is a set of securities that spans all market states. In this case we know that there are unique risk-neutral probabilities  $q_i$  for the market states.

If the prices are such that the project itself enters the optimal portfolio at zero level,  $U'_j$  is independent of the index  $j$ , and by the ratio theorem of the previous section, the risk-neutral probabilities  $q_{ij}$  are independent. Hence  $q_{ij}$  has the form  $q_{ij} = q_i^m p_j^n$ , where  $q_i^m$  is the risk-neutral probability for the market state, and  $p_j^n$  is the probability for the nonmarket state<sup>3</sup> (which is also the risk-neutral probability for that state).

---

<sup>3</sup> The independence argument applies even if there are more than two states in each part of the double tree.

Note that if there is no market component to a project, the project price (or value) is determined by its ordinary probabilities; that is, as the discounted expected value of its cash flows. At the other extreme, if the project has no private component, its price is determined by the risk-neutral market probabilities; that is, as the discounted risk-neutral expected value of its cash flows.

Here is a comprehensive example illustrating how these ideas can be used to evaluate a complex project. This example incorporates many of the concepts of this book and is worthy of careful study as an integrated review.

**Example 19.5 (Rapido: a rapidly declining oil well)** You are considering the possibility of investing in an oil well venture. If successful, the well life is likely to be about 25 years. The geological formations and other data indicate that this might be a favorable site. Before any initial drilling, the best estimate of the initial flow from the well if it is drilled is expressed as a list of possibilities and their probabilities, as shown in Table 19.1. We shall take a period length of 5 years in our analysis (to keep the problem size small enough to fit across a page). There are five possible levels of oil flows for the first 5 years of operation, which are shown in the table.

The initial drilling cost is \$220,000. After drilling, the initial flow can be estimated quite accurately, and a decision is then made as to whether to complete the well, making it ready for production. The completion cost is \$500,000. If the well is completed, the oil flow will decline as the reservoir is depleted. This decline can be expressed as a random chance that at the end of each 5-year period the flow will fall to the next lower category with a probability of 30%. This is a very rapid rate of decline for an oil well, and hence the name “Rapido.”

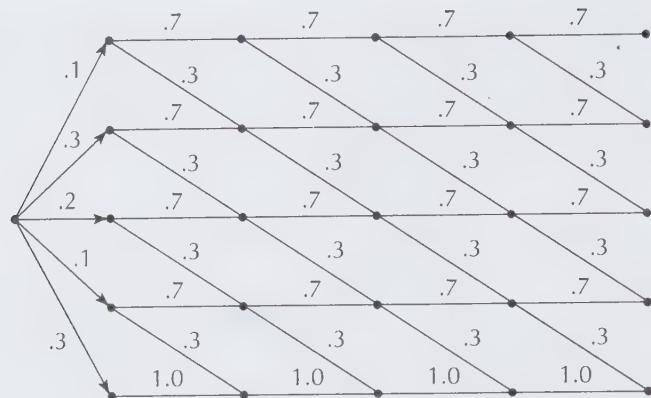
If the well is operated, the 5-year operating cost is \$400,000 of fixed cost plus \$5 per barrel in variable cost. All oil pumped from the well can be sold at the market price for crude oil, which is currently \$16 per barrel. We wish to find a fair price for this oil well venture, which has market risk associated with the future price of oil and technical (private) risk associated with the uncertainty of oil production.

The technical uncertainty regarding production possibilities is summarized by the lattice shown in Figure 19.7.

Next we must specify the market structure. For simplicity, we assume that the interest rate is constant at 7% per year or, equivalently, 40% for each 5-year period. It remains to specify the relevant aspects of the oil market. For this purpose we would first like an estimate of the volatility of oil prices. Such an estimate can be derived from a history of oil spot prices, but it is also possible to estimate the volatility directly from a single day’s record of option prices. There are no options for spot oil, but we

**TABLE 19.1**  
**OIL PRODUCTION POSSIBILITIES FOR THE INITIAL 5 YEARS**  
**OF OPERATION (IN THOUSANDS OF BARRELS PER 5-YEAR**  
**PERIOD).**

Oil produced	0	20	40	60	100
Probability	.3	.1	.2	.3	.1



**FIGURE 19.7 Technology of an oil well.** There are five possible levels of initial flow, which correspond to the five nodes that are successors to the initial node. The specific successor will be determined by the results of drilling. Then after each subsequent 5-year period, the oil flow either remains the same (with probability .7) or decreases one level (with probability .3).

OIL								METALS AND PETROLEUM									
CRUDE OIL (NYM) 1,000 BBLS.; \$ PER BBL.								CRUDE OIL, Light Sweet (NYM) 1,000 bbls.; \$ per bbl.									
Strike	Calls-Settle	Puts-Settle	Open	High	Low	Settle	Chg	Open	Lifetime	High	Low	Open					
Price	Jun	Jly	Aug	Jun	Jly	Aug		Jun	16.84	17.30	16.78	17.29	+	.43	21.35	14.02	124,032
1600	1.33	1.21	1.29	0.04	0.21	0.42		July	na	17.02	16.60	17.00	+	.35	20.78	14.15	73,360
1650	0.86	0.86	0.97	0.10	0.36	0.59		Aug	na	16.90	16.56	16.88	+	.30	20.78	14.35	34,123
1700	0.51	0.58	0.70	0.22	0.60	0.82		Sept	na	16.81	16.57	16.83	+	.27	20.78	14.50	28,809
1750	0.25	0.40	0.50	0.46	0.90	.....		Dec	16.57	16.80	16.55	16.80	+	.23	21.25	14.93	28,690
1800	0.11	0.25	0.34	0.82	1.25	.....		Jun	16.90	16.87	16.87	16.96	+	.18	21.21	15.73	17,396
1850	0.04	0.16	0.24	.....	.....	.....		Dec	.....	.....	.....	17.18	+	.16	20.80	16.50	10,793
Est vol	3,794	Wed	18,173	calls	8,785	puts		Jan	17.40	17.58	17.40	17.43	+	.14	20.26	17.22	14,698
Op int	Wed	211,586	calls	170,881	puts			Dec	.....	.....	.....	17.73	+	.13	20.40	17.53	19,072

**FIGURE 19.8 Quotes of oil future options and oil futures.** Volatility can be estimated option prices. Risk-neutral probabilities can be determined directly from futures market prices.

Sources: *Wall Street Journal*.

can use options on oil futures as a good substitute. A listing of these options is shown in the left table of Figure 19.8.

If we study the call options for August with strike prices of 1600 and 1700, we can use the Black–Scholes formula to solve for the implied volatility and the implied current futures price. This leads to an estimate of  $\sigma = 21\%$  (see Exercise 7), and we may assume that this is also the volatility of spot oil. If we use the standard binomial lattice approximation, we then set the up factor for oil at  $u = e^{\sigma \sqrt{\Delta t}} = e^{21\sqrt{5}} = 1.60$  and the down factor  $d = 1/1.60 = .625$ . (It is a great stretch of imagination to consider

$\Delta t = 5$  as “small”; however, we are treating this as a prototype model. A more complete model would use a smaller  $\Delta t$ .)

Now, usually, the next step would be to calculate the risk-neutral probabilities for this lattice using the formula  $q_u = (R - d)/(u - d)$ , giving  $q_u = .80$ , but this is *not* appropriate here. Oil has a significant storage cost; hence replication using oil would require paying storage costs. This will change the formula for risk-neutral probabilities. (See Section 15.11.) In fact, oil is generally not held as an investment, even though oil storage is possible, because the expected rate of return for doing so is not high enough to overcome the high storage costs. This tightness of the oil market is verified by the right side of Figure 19.8, which shows that the prices of oil futures contracts do not increase even as fast as the compounding of interest, as they would if markets were not tight. (See Section 12.3.) Indeed, we note that increasing the settlement date by  $2\frac{1}{2}$  years only increases the futures price by a factor of  $17.73/17.29 = 1.025$ . This is equivalent to about 1% per year.

We can, however, use the futures price information to determine appropriate risk-neutral probabilities. Given a spot price of  $S$ , next period the price will be either  $S_u$  or  $S_d$  according to our model. The current futures price for a contract that expires in 5 years will be about  $F = 1.05S$ . Since the current value of a futures contract is zero, and the payoff in 5 years will be either  $S_u - F$  or  $S_d - F$ , we must have

$$0 = q_u S(1.6 - 1.05) + q_d S(.62 - 1.05).$$

This yields

$$q_u = .44, \quad q_d = .56.$$

These are the values that we can use for the risk-neutral probabilities for oil price states.

We are now ready to carry out the backward recursion to determine the zero-level price of the oil venture. At the final period, from  $t = 20$  to  $t = 25$ , there are 25 possible states, corresponding to five oil flow components and five oil price components at that time. We think of these as being laid out in a 5 by 5 array. At the previous period there are the same five oil flow components and four oil price components, forming a 5 by 4 rectangle. This pattern progresses backward to period zero, just after completion of the well, where there is a 5 by 1 rectangle of states. Then, also at year 0, but before initial drilling, there is only a single node.

All of this is shown in Figure 19.9. To construct this figure the possible oil prices were first generated with a binomial lattice in the usual fashion, and these prices were laid out across the top row of the array according to the year in which they may occur. The possible flows were laid out down the last column of the array.

The backward calculation is a straightforward discounted expectation of cash flow and value. We assume for simplicity that all cash flow in a 5-year period occurs at the beginning of that period. Note that the final array consists only of profits from production in the last period. Earlier periods add current profit to a discounted risk-neutral expected value of the next period’s value. For example, the top right-hand

	$t = 0$		$t = 5$		$t = 10$			$t = 15$				$t = 20$				
Price	16	10	25.6	6.25	16	41	3.91	10	25.6	65.5	2.44	6.25	16	41	105	Flow
	1,938	517	3,994	67	1,523	6,713	0	279	2,756	9,395	0	0	700	3,196	9,586	100
	860	167	2,061	14.2	651	3,735	0	61	1,398	5,418	0	0	260	1,758	5,591	60
	288	46.9	1,000	1.94	203	2,085	0	8.8	694	3,292	0	0	40	1,038	3,594	40
Total	34.8	3.98	153	0	18.1	618	0	0	82.2	1,251	0	0	0	319	1,597	20
31.7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

**FIGURE 19.9 Rapido oil well evaluation.** The possible oil prices shown in the second row were generated by a binomial lattice, so the number of entries increases by one each period. There are five oil-flow possibilities each period. Backward evaluation is straightforward, once the proper risk-neutral probabilities are determined.

corner element in the array at  $t = 15$  is

$$\begin{aligned}
 v &= \text{flow} \times \text{oil price} - \text{cost} + \frac{1}{R} (\text{risk-neutral value of next period}) \\
 &= 100 \times 65.5 - 400 - 5 \times 100 + \frac{1}{1.4} (.44 \times .7 \times 9586 + .44 \times .3 \times 5591 \\
 &\quad + .56 \times .7 \times 3196 + .56 \times .3 \times 1758) \\
 &= 9395 \text{ (accounting for rounding errors).}
 \end{aligned}$$

The overall zero-level price accounts for the option to either complete the well or not. The zero-level price is \$31,700. Note how this rather complex problem is solved by a simple spreadsheet analysis—an analysis which, however, embodies a good deal of theory.

## 19.8 Buying Price Analysis

Frequently project opportunities arise in which investment must be either at a fixed positive level or at zero level, with nothing in between. An example is the opportunity to participate in a joint venture where each participant must subscribe to a fixed fraction of the project. Another is the prospect of taking on a project alone, such as the purchase of investment real estate. In such situations the zero-level price may not be the appropriate value, since the cash outlay required may represent a significant portion of one's investment capital.

A better concept of value in such situations is the **buying price**. The buying price is defined as the price that the investor would be willing to pay for participation in the project at the specified level. This price  $v_0$  is best understood in terms of expected utility. We first calculate the expected utility that would be achieved without participation in the project. Then we calculate the expected utility that would be achieved with participation, including an additional initial payment of an amount  $v_0$ . The value of  $v_0$  that makes these two expected utility values equal is the buying price. In other words, if  $v_0$  is the price to be paid for the project, the investor is indifferent between having the project or not. This price is different than the zero-level price,

which makes the investor indifferent between no participation and participation at a very small level.

## Certainty Equivalent and Exponential Utility

The buying price of a project can be computed easily if it is assumed that the investor's utility function is of exponential form,  $U(x) = -e^{-ax}$  for some  $a > 0$ . The computing procedure uses certainty equivalents rather than expected values.

Let us briefly review the certainty equivalent concept. Suppose that an investor has a utility function  $U$ . Suppose that  $X$  is a random variable describing the investor's wealth at the terminal point. Then the expected utility of this wealth is  $E[U(X)]$ . The certainty equivalent is the (nonrandom) amount  $\bar{x}$  such that  $U(\bar{x}) = E[U(X)]$ . We often write  $CE(X)$  for the certainty equivalent of  $X$ .

As a specific case suppose that  $U(X) = -e^{-ax}$  and suppose that the random variable  $X$  has two possible outcomes  $X_1$  and  $X_2$  occurring with probabilities  $p_1$  and  $p_2$ , respectively. The expected utility is

$$E[U(X)] = p_1 U(X_1) + p_2 U(X_2) = -p_1 e^{-aX_1} - p_2 e^{-aX_2}.$$

To find the certainty equivalent  $\bar{x}$  we solve

$$e^{-a\bar{x}} = p_1 e^{-aX_1} + p_2 e^{-aX_2}.$$

Taking the logarithm of both sides, we obtain

$$CE(X) = \bar{x} = -\frac{1}{a} \ln\{p_1 e^{-aX_1} + p_2 e^{-aX_2}\}. \quad (19.11)$$

This may look complicated, but it has a very special and important property.

The special property of this form is that if a constant, say  $\Delta$ , is added to a random variable, the certainty equivalent increases by this same constant. This property is often referred to as the **delta property**. Formally,

$$CE(X + \Delta) = CE(X) + \Delta$$

for any random variable  $X$  and any constant  $\Delta$ . This property can be checked easily for the two-outcome case by referencing equation (19.11).

Here is a general proof for exponential utility. We have

$$E(e^{-aX}) = e^{-aCE(X)}.$$

Therefore,

$$E[e^{-a(X+\Delta)}] = e^{-a\Delta} E(e^{-aX}) = e^{-a\Delta} e^{-aCE(X)} = e^{-a[CE(X)+\Delta]}.$$

This says that

$$\text{CE}(X + \Delta) = \text{CE}(X) + \Delta.$$

This delta property only holds for utility functions that are exponential or linear.

**Delta property** A utility function is linear or exponential if and only if for all random variables  $X$  and all constants  $\Delta$ , the certainty equivalent satisfies

$$\text{CE}(X + \Delta) = \text{CE}(X) + \Delta.$$

## Sequential Calculation of CE

Consider a one-period project having an initial known cash flow  $c_0$  followed at the end of the period by a random cash flow that takes one of the values  $c_1$  or  $c_2$  with probabilities  $p_1$  and  $p_2$ , respectively. There is also a risk-free asset with return  $R$ . This project is illustrated in Figure 19.10.

Assume that the investor has initial wealth  $X_0$  and uses an exponential utility on final wealth. Risk-free borrowing or lending is used to transfer any cash flow at the initial time to a cash flow at the final time. If the project is not taken, then the final utility value will be  $U(RX_0)$  since the initial wealth is transformed by the risk-free return.

If the project is taken at a price  $v_0$ , the expected utility of final wealth will be

$$p_1 U\{[c_1 + R(X_0 + c_0 - v_0)]\} + p_2 U\{[c_2 + R(X_0 + c_0 - v_0)]\}.$$

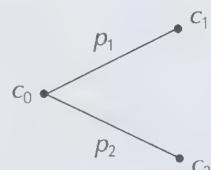
When the price  $v_0$  is set correctly, the expected utility with the project will equal the value without the project, namely,  $U(RX_0)$ . Setting the certainty equivalents of these two equal to each other, we obtain<sup>4</sup>

$$\text{CE}[c_1 + R(x_0 + c_0 - v_0), c_2 + R(x_0 + c_0 - v_0)] = RX_0.$$

Note that both terms on the left contain  $R(x_0 + c_0 - v_0)$ . This is a constant, and by the delta property it can be taken out of the CE expression. We therefore obtain

$$\text{CE}[c_1, c_2] + R(X_0 + c_0 - v_0) = RX_0.$$

**FIGURE 19.10 Simple project.** This project has initial cash flow  $c_0$ , followed at the end of the period by a cash flow of value either  $c_1$  or  $c_2$ .



<sup>4</sup> As a shorthand notation, if  $c_1$  and  $c_2$  are cash flows in two final states, we write  $\text{CE}[c_1, c_2]$  for the corresponding certainty equivalent.

Solving for  $v_0$ , we obtain an expression for the buying price,

$$v_0 = c_0 + \frac{1}{R} \text{CE}[c_1, c_2]. \quad (19.12)$$

Note that this equation looks just like a net present-value formula. The certainty equivalent is used to summarize the cash flow at the end of the period.

## Multiperiod Case

The preceding technique extends to cash flow processes defined over several periods, but the risk aversion coefficient of the utility function must be adjusted each period. Specifically, the risk aversion coefficient used to evaluate the certainty equivalent at time  $t$  must be  $aR^{T-t}$  instead of the original  $a$ . This reflects the fact that the effective utility function for money  $X$  received at time  $t$  is  $U(R^{T-t}X)$  rather than  $U(X)$  because  $X$  will be transformed to  $R^{T-t}X$  at time  $T$ .

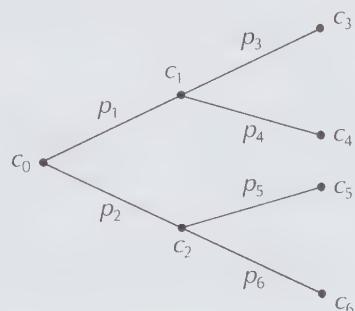
As an example of the full calculation, consider the two-period project shown in Figure 19.11. To evaluate this project we work backward in the usual fashion. First we calculate  $v_1$  at the node where  $c_1$  occurs by using the formula for the one-period case; namely,  $v_1 = c_1 + (1/R)\text{CE}_2[c_3, c_4]$ , where the subscript on CE denotes that the appropriate risk aversion coefficient at  $t = 2$  (which is  $a$ ) is used. Next  $v_2$  is computed at the  $c_2$  node in an analogous fashion as  $v_2 = c_2 + (1/R)\text{CE}_2[c_5, c_6]$ . Finally, we find

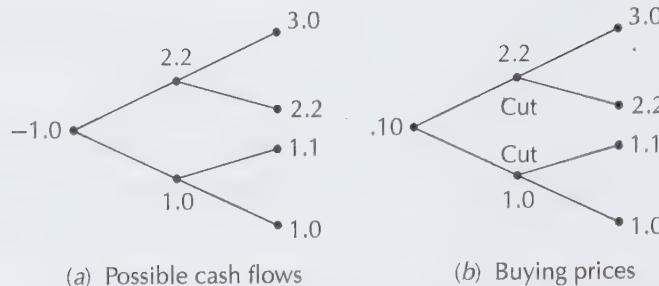
$$v_0 = c_0 + \frac{1}{R} \text{CE}_1[v_1, v_2]. \quad (19.13)$$

This final certainty equivalent is computed with the risk aversion coefficient magnified by one period of interest, and with the probabilities  $p_1$  and  $p_2$  for  $v_1$  and  $v_2$ , respectively.

**Example 19.6 (When to cut a tree)** Consider again the tree-cutting example treated in the last section, but this time suppose that we are planning to purchase this project ourselves. We must buy the full project or none of it. The project cash flow possibilities are shown in Figure 19.12(a). Recall that the figures at the intermediate

**FIGURE 19.11 Two-period project.** The buying price can be found by evaluating the certainty equivalent by a backward process.





**FIGURE 19.12** Buying price for tree farm. (a) Cash flow generated at a node if the trees are cut at that point. (b) Certainty equivalent at a node and best policy.

nodes are the cash flows that would be attained if the trees were cut there and the process terminated. Also, all arcs have probability .5.

Assume that our utility function is  $U(x) = -e^{-3x}$  and the interest rate is 10% per year, as in the earlier example. The first step is to calculate the certainty equivalent of the last two upper nodes. This certainty equivalent is

$$-\frac{1}{3} \ln [ .5e^{-3 \times 3.0} + .5e^{-3 \times 2.2}] = 2.4.$$

When discounted one period, this becomes 2.18. Since this is less than the 2.2 value that would be achieved by cutting the trees at that point, we decide to cut, and we assign the buying price of 2.2 to that node. The node below that also retains the value of 1.0, since it is clear that the discounted certainty equivalent of the lower last phase is less than 1.

Finally, we calculate the buying price at the first node. To calculate the certainty equivalent, we must change the risk aversion coefficient from  $a$  to  $aR$ , or in this case from 3 to 3.3. Accordingly, the proper utility function for this period is  $U(x) = -e^{-3.3x}$ . Hence the certainty equivalent of the middle two nodes is

$$-\frac{1}{3.3} \ln [ .5e^{-3.3 \times 2.2} + .5e^{-3.3 \times 1.0}] = 1.21.$$

Discounting this and accounting for the original cash flow, we find  $v_0 = .10$ . This is quite a bit lower than the zero-level price of .52 found in the last section. The price must be lower to induce us to purchase the entire project rather than just a small fraction of it.

## General Approach

Suppose now that states of the world can be factored into independent market and nonmarket components. A general state at time  $t$  is written as in the last section as  $(s_t^m, s_t^n)$ , corresponding to the market and nonmarket components. We also assume that the market portion of the system is complete; that is, there is a complete set of

assets that span all dimensions of the market. In that case we know that there are unique risk-neutral probabilities  $q_i$  for the market states.

We assume that the investor has an exponential utility function for final wealth. The project has cash flows specified at each node.

To find the buying price, we proceed recursively, starting at the final time. At the final time the buying price at any node is equal to the cash flow at that node. At any other (previous) node  $(s_t^m, s_t^n)$  of the backward process, two calculation steps are required. First, for each fixed market successor  $i$ , we compute the certainty equivalent with respect to the nonmarket components  $j$ . That is, we find the certainty equivalent  $\text{CE}_i$  such that  $U(R^{T-t}\text{CE}_i) = \sum_j p_j^n U(R^{T-t}v_{ij})$ , where  $v_{ij}$  is the buying price of the successor node  $ij$ . Then we find the new buying price from

$$v_t^{m,n} = c_t^{m,n} + \frac{1}{R} \sum_i q_i \text{CE}_i.$$

In other words, we use certainty equivalent calculation on the nonmarket component and risk-neutral pricing on the market component.

**Example 19.7 (Rapido oil well)** We can analyze the Rapido oil well using a certainty equivalent analysis. Only a few modifications to the earlier zero-level price analysis are required. We assume that a single investor is planning to finance the entire project. This investor has a utility function  $U(X) = -e^{-X/10,000}$ , where  $X$  is in thousands of dollars. This is realistic for an investor having a net worth of about \$10 million (20 years from now).

In order to find the buying price, we simply change the risk-neutral discounting formula to one that is a mixture of risk-neutral pricing of the market state (the oil price) and a certainty equivalent of the technical factors (the flow level). We must remember to update the effective utility function by the factor of  $R = 1.4$  in the exponent each period. The results are shown in Figure 19.13.

The final array, at  $t = 20$ , is identical to that of the earlier example, since that array contains final cash flows. The upper right-hand corner element of the array at

	$t = 0$		$t = 5$			$t = 10$			$t = 15$			$t = 20$				
Price	16	10	25.6	6.25	16	41	3.91	10	25.6	65.5	2.44	6.25	16	41	105	Flow
	1,900	512	3,926	66.8	1,514	6,624	0	278	2,748	9,331	0	0	700	3,196	9,586	100
	849	165	2,041	14.2	649	3,710	0	60.8	1,396	5,402	0	0	260	1,758	5,591	60
	281	46.7	985	1.94	201	2,063	0	8.79	692	3,276	0	0	40	1,038	3,594	40
Total	34.5	3.98	150	0	18.1	610	0	0	81.9	1,242	0	0	0	319	1,597	20
-5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

**FIGURE 19.13 Certainty equivalent analysis of Rapido oil well.** A complex problem is treated by a single spreadsheet model. Vertical pairs are combined by certainty equivalents, and horizontal pairs of these are combined by risk-neutral probabilities.

$t = 15$  is evaluated as

$$v = \text{flow} \times \text{oil price} - \text{cost}$$

$$+ \frac{1}{R} [q_u \text{CE}(9586, 5591) + q_d \text{CE}(3196, 1758)].$$

We have

$$\text{CE}(9586, 5591) = -10,000 \times \ln[.7e^{-.9586} + .3e^{-.5591}] = 8211$$

$$\text{CE}(3196, 1758) = -10,000 \times \ln[.7e^{-.3196} + .3e^{-.1758}] = 2742.$$

Then using  $q_u = .44$ , and  $q_d = .56$  from the earlier example, we obtain

$$v = 100 \times 65.5 - 400 - 5.100 + \frac{.44 \times 8211 + .56 \times 2742}{1.4}$$

$$= 9331.$$

Note that the initial buying price is negative, which indicates that the project is too big for this investor to take on alone. It is a good project, as shown by the zero-level analysis, but only when a smaller share is taken or a smaller risk aversion coefficient is used.

## 19.9 Pricing Axioms for Continuous Time

The Black–Scholes equation for pricing derivative securities in continuous time can be generalized to price assets that are not pure derivatives. The generalization uses projection pricing, defined by a generalized pricing operator and obeying a set of pricing axioms that extend those of Section 15.13. The end result includes a combination of the two most celebrated pricing methods: Black–Scholes and CAPM. The result is robust, in that every risk-averse investor will, in a specific sense, agree that the resulting price process is appropriate.

In the generalized version, we consider the **continuous-time environment** that includes several marketed assets and a variable  $x_e$  outside of this market that follows geometric Brownian motion. The entire environment of this variable and  $n$  marketed assets is described by the system

$$dx_e = \mu_e x_e dt + \sigma_e x_e dz_e$$

$$dx_i = \mu_i x_i dt + \sigma_i x_i dz_i$$

for  $i = 1, 2, \dots, n$ , where the  $x_i$ 's are prices of marketed assets. In addition, there is a risk-free asset in the market with rate  $r$ . The first equation describes  $x_e$ , which cannot be traded. The normalized market Wiener process inputs are correlated as  $E[dz_i dz_j] = \rho_{ij} dt$ . We also have  $E[dz_e dz_i] = \rho_{ei} dt$  for all  $i$ . We define  $\sigma_{em} = \rho_{em}\sigma_e\sigma_m$  and denote by  $M$  the space of marketed assets—basically all portfolios of the  $n + 1$  marketed assets described. We assume that no arbitrage possibility exists in  $M$ .

There is a terminal payoff function  $F(x_e, T)$  that is to be priced. We seek an overall value function  $V(x_e, t)$  that gives the value for all feasible  $x_e$  and  $t \in [0, T]$ . By definition, we have the boundary condition  $V(x_e, T) = F(x_e, T)$ .

For the example mentioned earlier, the variable  $x_e$  might represent the revenue of a firm and define the payoff of an option. Revenue is not traded, but it may be correlated with assets that are traded, such as the firm's stock price.

Other examples arise in the context of **real options** where progress of a project to time  $t$  is measured by a variable  $x_e$ , with final project payoff being determined by this variable.

In this type of environment, the pricing operator  $\mathbb{P}$  is generalized by projecting onto the market at each instant. This is done by changing Axiom 1 to

$$(1) \quad x = \mathbb{P}\{x + dx | M\},$$

which means that the argument of  $\mathbb{P}$  must be projected onto the marketed space. Essentially, if  $dx$  is outside of  $M$ , only the projected portion is used in pricing.

The other axioms remain the same; namely,

- (2) For a constant  $C$ ,  $\mathbb{P}\{C\} = C(1 - r dt)$
- (3)  $\mathbb{P}$  is linear
- (4)  $\mathbb{P}\{dt\} = dt$ . Terms of higher order in  $dt$  are ignored.

We shall apply these axioms to derive (relatively quickly) an appropriate pricing equation. Here is the final result:

**Theorem 19.1 Generalized Black–Scholes equation** Suppose that the variable  $x_e$  determines the terminal reward  $F(x_e, T)$ . Let  $V(x_e, t)$  be the projection price of this payoff at time  $t < T$  and variable value  $x_e(t)$ . Let  $x_m$  be a portfolio of marketed assets whose return is most correlated with  $dx_e$ . Then  $V(x_e, t)$  satisfies the equation

$$rV(x_e, t) = V_t(x_e, t) + V_{x_e}(x_e, t)x_e[\mu_e - \beta_{em}(\mu_m - r)] + \frac{1}{2}V_{x_ex_e}(x_e, t)x_e^2\sigma_e^2, \quad (19.14)$$

with  $V(x_e, T) = F(x_e, T)$  and  $\beta_{em} = \sigma_{em}/\sigma_m^2$ .

Notice the similarity of equation (19.14) to the standard Black–Scholes equation. Indeed, if the variable  $x_e$  is in the market, then  $x_m$  can be taken as  $x_e$  and  $\beta_{em}$  will be 1. The bracketed term on the right will reduce to  $r$ , corresponding to the term in the Black–Scholes equation. In the general case where  $x_e$  is not a market variable, and hence not a derivative, the equation will have a term that includes  $\mu_e$ .

The proof given here of this theorem is based on the pricing axioms.

**Proof:** From Ito's equation the value function  $V(x_e, t)$  satisfies

$$dV = [V_t + V_{x_e}\mu_{ex_e} + \frac{1}{2}V_{x_ex_e}x_e^2\sigma_e^2]dt + V_{x_e}\sigma_{ex_e}dz_e. \quad (19.15)$$

Let  $m$  be a marketed portfolio most correlated with  $x_e$ . We know that the projection of  $dz_e$  onto  $M$ , denoted  $\{dz_e | M\}$ , satisfies  $\{dz_e | M\} = ar dt + b dz_m$  for some

constants  $a$  and  $b$ . Minimizing  $E(dz_e - ar dt - bdz_m)^2$  shows that

$$\{dz_e|M\} = \rho_{em} dz_m. \quad (19.16)$$

Using equation (15.42) for the market price of risk, we have

$$\mathbb{P}\{dz_m\} = \frac{(r - \mu_m)}{\sigma_m} dt.$$

Substituting from equation (19.16), we get

$$\mathbb{P}\{dz_e|M\} = \frac{(r - \mu_m)}{\sigma_m} \rho_{em} dt.$$

With  $\sigma_{em} = \rho_{em} \sigma_e \sigma_m$  and  $\beta_{em} = \sigma_{em}/\sigma_m^2$ , we may use the alternate form

$$\mathbb{P}\{dz_e|M\} = \frac{\beta_{em}}{\sigma_e} (r - \mu_m) dt. \quad (19.17)$$

As in equation (15.38), we have<sup>5</sup>  $\mathbb{P}\{dV|M\} = rVdt$ .

Returning to equation (19.15), we may apply the operator  $\mathbb{P}$  (with projection) to both sides. From earlier, the left is  $rV(x_e, t)dt$ . We know  $\mathbb{P}\{dz_e|M\}$  from equation (19.17). Using that and canceling  $dt$  we obtain the extended Black–Scholes equation,

$$rV(x_e, t) = V_t(x_e, t) + V_{x_e}(x_e, t)x_e[\mu_e - \beta_{em}(\mu_m - r)] + \frac{1}{2}V_{x_ex_e}(x_e, t)x_e^2\sigma_e^2, \quad (19.18)$$

with  $V(x_e, T) = F(x_e, T)$ . ■

From our study in earlier chapters we know that projection pricing is popular, although there may be important alternatives. However, in the present case, a strong argument can be made that the value function  $V(x_e, t)$  is a most natural choice. In particular, if the price of the variable  $x_e$  is preassigned to be  $V(x_e, t)$  and a risk-averse investor is allowed to construct a portfolio process consisting of market assets and the variable  $x_e$ , the investor will never find it advantageous to include  $x_e$ . The price function  $V(x_e, t)$  is a universal zero-level price.

To see this easily, note that  $dx_e = dx_{ep} + dz_{ef}$  where  $dx_{ep}$  is the projection of  $dx_e$  onto the market, and  $dz_{ef}$  is uncorrelated with the market. Since  $dx_{ep}$  is a marketed variable, its zero-level price is its market price, which is  $V(x_e, t)$  by projection pricing. The zero-level price of  $dz_{ef}$  is zero since it is independent of the market and has zero expected value. Hence, by linearity, the zero-level price of  $dx_e$  is  $V(x_e, t)$ . It is in fact a universal zero-level price. (See Section 11.9.)

A nice way to compare this with the standard Black–Scholes result as follows: The pricing of a derivative by the standard Black–Scholes framework renders the derivative *redundant* because it can be replicated exactly. In the generalized case, the pricing of a derivative renders the derivative *irrelevant* because no risk-averse investor will benefit by including it in an optimal portfolio.

---

<sup>5</sup>  $\mathbb{P}\{dV|M\} = \mathbb{P}\{V + dV|M\} - \mathbb{P}\{V|M\} = V - (1 - r dt)V = rVdt$ .

## Option Formula

In the case of an option based on a nonmarketed variable, a modification of the standard Black–Scholes formula applies. Specifically, for an asset  $x_e$  and a terminal value of  $\max(0, x_e(T) - K)$ , we denote  $\omega_e = \mu_e - \beta_{e,m}(\mu_m - r)$ . Then (with  $x_e = S$ )

$$V = Se^{(\omega_e - r)T}N(d_1) - Ke^{-rT}N(d_2),$$

where

$$d_1 = \frac{\ln(S/K) + (\omega_e + \frac{1}{2}\sigma_e^2)T}{\sigma\sqrt{T}} \quad (19.19)$$

$$d_2 = d_1 - \sigma\sqrt{T}. \quad (19.20)$$

## Risk-Neutral Form

From the generalized Black–Scholes equation we can deduce that the risk-neutral version of the  $x_e$  process is

$$dx_e = \omega_e x_e dt + \sigma_e x_e dz_e \quad (19.21)$$

where  $\omega_e = \mu_e - \beta_{em}(\mu_e - r)$ . With this form, we can find the value of a terminal function  $F(x_e, (T))$  as

$$p = e^{-rT}\hat{E}[F(x_e(T))]. \quad (19.22)$$

where  $\hat{E}$  denotes expectation at time 0 with respect to the risk-neutral process. This formulation can provide the basis of determining value by simulation.

## Alternative Forms

The pricing formula in Theorem 19.1 uses a marketed asset most correlated with the variable  $x_e$  as a basis for pricing. This method parallels the correlation pricing method of Chapter 7. However, similar to the situation in Chapter 7, certain other marketed portfolios can be used in place of a most-correlated portfolio and the result will be the same. For example we may define a (generalized) Markowitz portfolio  $x_M$  as the marketed portfolio of risky assets that maximizes  $(\mu_M - r)/\sigma_M$ . When it exists, this portfolio can be used in place of a most-correlated asset. This has the advantage that it will price any variable-based derivative, not just the one based on the given variable  $x_e$ . In addition, a portfolio analogous to the minimum-norm portfolio can serve this same purpose.

On the other hand, for purpose of approximating a derivative of  $x_e$  with market assets or for hedging such a derivative if it is owned, a most-correlated market asset is likely to be most useful.

## 19.10 Summary

Evaluation of an investment opportunity is an extension of the concept of present value. This value is based on the relationship of the cash flow of the opportunity to the cash flows and values of marketed assets.

The construction of a graph to represent a group of assets can be a challenge. One approach is to start with binomial lattice representations of each asset separately, and combine them into a double, triple, or multilattice in such a way as to capture the covariance structure of the assets. The details are easily worked out for a double lattice. Once the risk-neutral probabilities are determined, the price of a security can be found by the backward process of discounted risk-neutral valuation.

Private uncertainty is treated differently from market uncertainty because there are no associated market prices. Usually this means that the actual private probabilities should be used just like risk-neutral probabilities to determine the zero-level price of an asset.

The buying price of a project or asset is the price that an investor would pay to accept the project or asset in full (or a specified portion of it). This price depends on the investor's utility function and is usually lower than the zero-level price. If the utility function for final wealth is exponential, a backward evaluation process can be used to find the buying price. This procedure uses certainty equivalents to evaluate private uncertainty and risk-neutral prices to evaluate market uncertainties. This is because the private uncertainty cannot be hedged, but the market uncertainty can.

In the case of continuous time, one approach is to project locally onto the marketed assets in a manner similar to the CAPM method. This leads to a partial differential equation that generalizes the Black–Scholes equation.

## Exercises

1. (A state tree) A certain underlying state graph is a tree where each node has three successor nodes, indexed  $a$ ,  $b$ ,  $c$ . There are two assets defined on this tree which pay no dividends except at the terminal time  $T$ . At a certain period it is known that the prices of the two assets are multiplied by factors, depending on the successor node. These factors are shown in Table 19.2.

- (a) Is there a short-term riskless asset for this period?
- (b) Is it possible to construct an arbitrage?

**TABLE 19.2**

	<i>a</i>	<i>b</i>	<i>c</i>
Security 1	1	1.2	1.0
Security 2	2	1.2	1.3

2. (Node separation) Consider a short rate binomial lattice where the risk-free rate at  $t = 0$  is 10%. At  $t = 1$  the rate is either 10% (for the upper node) or 0% (for the lower node). Trace out the growth of \$1 invested risk free at  $t = 0$  and rolled over at  $t = 1$  for one more period. The values obtained at  $t = 1$  and  $t = 2$  correspond to  $R_{01}$  and  $R_{02}$ . Show that these factors cannot be represented on a binomial lattice, but rather a full tree is required. Draw the tree.

- 3. (Bond valuation)** Assuming the short rate process of Exercise 2 and risk-neutral probabilities of .5, consider a zero-coupon bond that pays \$1 at time  $t = 2$ . Find the value at time  $t = 0$  of this bond in two ways:
- Using the short rate lattice and equation (19.1).
  - Using the tree for  $R_{0s}$  and equation (19.2).
- 4. (Optimal option valuation  $\oplus$ )** Find the values of the 5-month call option of Example 19.1 using the same trinomial lattice used in that example but employing the utility function  $U(x) = \sqrt{x}$ . What is  $\alpha$ ?
- 5. (Gold correlation)** Suppose that in the double stochastic Simplico gold mine example the real probability of an up move in gold is .6 and the real probability of an up move in the short rate is .7. Suppose also that gold price and short rate fluctuations have a correlation coefficient of  $-.4$ . Find the appropriate  $q_{ij}$ 's.
- 6. (Complexico mine  $\oplus$ )** Use the information about the Complexico mine of Example 14.8, Chapter 14, but assume that gold prices and interest rates are governed by the models of Example 19.3. Find the value of the Complexico lease.
- 7. (Simultaneous solution)** Calculate the volatility and the current price of oil futures implied by the call 1600 August and the call 1700 August of Figure 19.8 by using the Black–Scholes formula with  $T = .25$ .
- 8. (Default risk  $\oplus$ )** A company issues a 10% coupon bond that matures in 5 years. However, this company is in trouble, and it is estimated that each year there is a probability of .1 that it will default that year. (Once it defaults, no further coupons or principal are paid.) What is the value of the bond?
  - Assume the term structure of interest is flat at 10%.
  - Assume that the short rate is currently 10% and the short rate is multiplied by either 1.2 or .9 each year with risk-neutral probabilities of .5. Default risk is independent of the interest rate.
- 9. (Automobile choice)** Mr. Smith wants to buy a car and is deciding between brands A and B. Car A costs \$20,000, and Mr. Smith estimates that at the rate he drives he will sell it after 2 years and buy another of the same type for the same price. The resale price will be either \$10,000 or \$5,000, each with probability .5, at the end of each 2-year period. Car B costs \$35,000 and will be sold after 4 years with an estimated resale price of either \$12,000 or \$8,000, each with probability .5. The yearly maintenance costs of the two cars are constant each year and identical for the two cars. Mr. Smith has an exponential utility function with risk aversion coefficient of about  $a = 1/\$1,000$  now. Real interest is constant at 5%. Which car should he decide is better from an economic perspective over a 4-year period, and what is the certainty equivalent of the difference?
- 10. (V example)** Consider a continuous-time environment, with  $e$  as a variable outside of the market.
  - Suppose the final payoff is  $V(x_e, T) = x_e(T)$ . Find  $V(x_e, t)$ .
  - Find a put–call relation that holds for options on a nontraded asset  $e$ .
- 11. (Gavin's final)** Mr. Jones was considering a new grapefruit venture that would generate a random sequence of yearly cash flows. He asked his son, Gavin, “People tell me I should use a cost of capital figure to discount the stream. They say it's based on the CAPM. Have you given up on that? I haven't heard you talk about it for awhile.”
- Gavin replied, “Special conditions are required to justify it for more than one period. We had a complicated final exam question on it.”

Consider a two-year model. The risk-free rate for each is  $r$ . The (random) rates of return for the Markowitz portfolio in the two years are  $r_1$  and  $r_2$ , respectively, and they are independent. There is a single random cash flow  $x_2$  at the end of the second year. Denote by  $x_{2|0}$  and  $x_{2|1}$  the random variable  $x_2$  given the information at times zero and one, respectively, and let  $E_0$  and  $E_1$  denote expectation at times zero and one. Likewise let  $V_0$  and  $V_1$  denote the value at time zero and one, respectively, of receiving  $x_2$  at time 2. Assume that  $E_0[E_1[x_{2|1}]] = E_0[x_{2|0}]$  and that  $\text{cov}[x_{2|1}/V_1, r_2]$  is independent of the information received at time one. Show that the value at time zero of receiving  $x_2$  at time 2 is

$$V_0 = \frac{E_0[x_{2|0}]}{[1 + r + \beta_1(\bar{r}_1 - r)][1 + r + \beta_2(\bar{r}_2 - r)]},$$

where

$$\beta_1 = \text{cov}[V_1/V_0, r_1]/\sigma_{r_1}^2.$$

Find  $V_1$  and  $\beta_2$ .

- 12.** (Soft option) An option based on a variable that is not traded is called a **real option** or sometimes a **soft option**. Find the projection price of the soft option with the following parameters and compare with the option of Example 15.2.

$$K = 60, \quad S = 62, \quad r = 10\%, \quad \sigma = 20\%, \quad T = 5 \text{ months}, \\ \mu_e = 8\%, \quad \sigma_m = 15\%, \quad \mu_m = 14\%, \quad \rho = 0.7$$

## References

The overall structure of multiperiod investments is presented comprehensively in Duffie [1]. Construction of multiperiod lattices has been approached in several ways. See for example [2–3]. For the theory here see [4]. The buying price analysis is adapted from Smith and Nau [5]. Also see [6]. For a general theory of pricing see Holton [7]. For pricing an option as the marginal (zero-level) price see Davis [8]. This price was shown to be universal in [9]. For projection pricing and the generalized Black–Scholes equation, see [10]. For the discrete-time version, see [11].

1. Duffie, D. (1996), *Dynamic Asset Pricing Theory*, 2nd ed., Princeton University Press, NJ.
2. Boyle, P. P., J. Evnine, and S. Gibbs (1989), “Numerical Evaluation of Multivariate Contingent Claims,” *Review of Financial Studies*, **2**, 241–250.
3. He, H. (1990), “Convergence from Discrete- to Continuous-Time Contingent Claims Prices,” *Review of Financial Studies*, **3**, 523–546.
4. Luenberger, D. G. (1996), “Double Trees for Investment Analysis,” presented at the Conference on Computational Economics and Finance, Geneva, June.
5. Smith, J. E., and R. F. Nau (1995), “Valuing Risky Projects: Option Pricing Theory and Decision Analysis,” *Management Science*, **41**, no. 5, 795–816.
6. Carmona, R. (ed.) (2009) *Indifference Pricing*, Princeton University Press, Princeton, New Jersey.
7. Holton, H. H. (1997). *Asset Valuation and Optimal Portfolio Choice in Incomplete Markets*, Ph.D dissertation, Department of Engineering–Economic Systems, Stanford University, Stanford, CA.
8. Davis, M. H. A. (1997) “Option Pricing in Incomplete Markets,” In M. A. H. Dempster and S. R. Pliska, (editors), *Mathematics of Derivative Securities*, 216–227. Cambridge University Press.
9. Karatzas, I. and S. G. Lou, (1996) “On the pricing of contingent claims under constraints”, *Annals of Applied Probability*, **6**, 321–369.
10. Luenberger, D. G. (2004), “Pricing a Nontradeable Asset and Its Derivatives,” *Journal of Optimization Theory and Applications*, **121**, 465–487.
11. Luenberger, D. G. (2011), “Pricing Dynamic Binary Variables and Their Derivatives,” *Quantitative Finance*, (2012) vol 12, 451–464.

# APPENDIX A

## BASIC PROBABILITY THEORY

### A.1 General Concepts

As discussed in Chapter 6, a random variable  $x$  is described by its probability structure. If  $x$  can take on only a finite number of values, say,  $x_1, x_2, \dots, x_m$ , then the density function gives the probability for each of those outcome values. We frequently use  $P$  to denote probability. Thus

$$P(x_i) = \text{prob}(x_i);$$

that is,  $P(x_i)$  is the probability that  $x$  takes on the value  $x_i$ . We always have  $P(x_i) \geq 0$  for all  $x_i$ . Also,  $\sum_i P(x_i) = 1$ .

If the random variable  $x$  can take on a continuum of values, such as all real numbers, then the probability density function  $p(\xi)$  is also defined for all these values. The interpretation in this case is, roughly, that

$$p(\xi)d\xi = \text{prob}(\xi \leq x \leq \xi + d\xi) = P(\xi \leq x \leq \xi + d\xi).$$

The **probability distribution** of the random variable  $x$  is the function  $F(\xi)$  defined as

$$F(\xi) = P(x \leq \xi).$$

It follows that  $F(-\infty) = 0$  and  $F(\infty) = 1$ . In the case of a continuum of values, if  $F$  is differentiable at  $\xi$ , then  $dF(\xi)/d\xi = p(\xi)$ .

Two random variables  $x$  and  $y$  are described by their **joint probability density** or **joint probability distribution**. The joint distribution is the function  $F$  defined as

$$F(\xi, \eta) = P(x \leq \xi, y \leq \eta).$$

The joint density is defined in terms of derivatives, or if there are only a finite number of possible outcomes, the joint density at a pair  $x_i, y_i$  is  $P(x_i, y_i)$  equal to the probability of that pair occurring. In general,  $n$  random variables are defined by their joint probability distribution defined with respect to  $n$  variables.

From a joint distribution the distribution of any one of the random variables can be easily recovered. For example, given the distribution  $F(\xi, \eta)$  of  $x$  and  $y$ , the distribution of  $x$  is

$$F_x(\xi) = F(\xi, \infty).$$

The random variables  $x$  and  $y$  are **independent** if the density function factors into the form

$$p(\xi, \eta) = p_x(\xi)p_y(\eta).$$

This is the case for the pair of random variables defined as the outcomes on two fair tosses of a die. For example, the probability of obtaining the pair (3, 5), in order, is  $\frac{1}{6} \times \frac{1}{6}$ .

The **expected value** of a random variable  $x$  with density function  $p$  is

$$E(x) = \int_{-\infty}^{\infty} \xi p(\xi) d\xi.$$

If  $E(x)$  is denoted by  $\bar{x}$ , the **variance** of  $x$  is

$$\text{var}(x) = \int_{-\infty}^{\infty} (\xi - \bar{x})^2 p(\xi) d\xi.$$

Likewise, the **covariance** of  $x$  and  $y$  is

$$\text{cov}(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\xi - \bar{x})(\eta - \bar{y}) p(\xi, \eta) d\xi d\eta.$$

It is easy to show that if  $x$  and  $y$  are independent, then they have zero covariance.

## A.2 Normal Random Variables

A random variable  $x$  is said to be **normal** or **Gaussian** if its probability density function is of the form

$$p(\xi) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(\xi-\mu)^2}.$$

In this case the expected value of  $x$  is  $\bar{x} = \mu$  and the variance of  $x$  is  $\sigma^2$ . This density function is the characteristic “bell-shaped” curve, illustrated in Figure A.1.

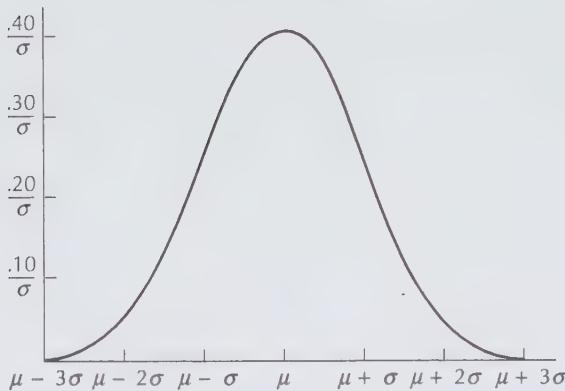
A normal random variable is **normalized** or **standard** if  $\bar{x} = 0$  and  $\sigma^2 = 1$ . Thus a standard normal random variable has the density function (written in terms of the variable  $x$ )

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}.$$

The corresponding standard distribution is denoted by  $N$  and given by the expression

$$N(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}\xi^2} d\xi.$$

There is no analytic expression for  $N(x)$ , but because of its importance, tables of its values and analytic approximations are available.



**FIGURE A.1 Normal distribution.** The expected value is  $\mu$  and the variance is  $\sigma^2$ .

To work with more than one normal random variable it is convenient to use matrix notation. We let  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  be a vector of  $n$  random variables. The expected value of this vector is the vector  $\bar{\mathbf{x}}$ , whose components are the expected values of the components of  $\mathbf{x}$ . The **covariance matrix** associated with  $\mathbf{x}$  is the  $n \times n$  matrix  $\mathbf{Q}$  with components  $[\mathbf{Q}]_{ij} = \text{cov}(x_i, x_j)$ . If  $\mathbf{x}$  is regarded as a column vector and  $\mathbf{x}^T$  is the corresponding row vector, then  $\mathbf{Q}$  can be expressed as

$$\mathbf{Q} = E[(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T].$$

If the  $n$  variables are jointly normal, the distribution of  $\mathbf{x}$  is

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{Q}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{Q}^{-1} (\mathbf{x} - \bar{\mathbf{x}})}.$$

If two jointly normal random variables are uncorrelated, then it is easy to see that the joint density function factors into a product of densities for the two separate variables. Hence if two jointly normal random variables are uncorrelated, they are independent.

A most important property of jointly normal random variables is the summation property. Specifically, if  $x$  and  $y$  are jointly normal, then all random variables of the form  $\alpha x + \beta y$ , where  $\alpha$  and  $\beta$  are constants, are also normal. This result is easily extended to higher order sums. In fact if  $\mathbf{x}$  is a column vector of jointly normal random variables and  $\mathbf{T}$  is an  $m \times n$  matrix, then the vector  $\mathbf{T}\mathbf{x}$  is an  $m$ -dimensional vector of jointly normal random variables.

## A.3 Lognormal Random Variables

A random variable  $z$  is lognormal if the random variable  $\ln z$  is normal. Equivalently, if  $x$  is normal, then  $z = e^x$  is lognormal. In concrete terms this means that the density

function for  $z$  has the form

$$p(\zeta) = \frac{1}{\sqrt{2\pi}\sigma\zeta} e^{-\frac{1}{2\sigma^2}(\ln\zeta - v)^2}.$$

We have the following values:

$$\mathbb{E}(z) = e^{(v+\sigma^2/2)} \quad (\text{A.1})$$

$$\mathbb{E}(\ln z) = v \quad (\text{A.2})$$

$$\text{var}(z) = e^{(2v+\sigma^2)}(e^{\sigma^2} - 1) \quad (\text{A.3})$$

$$\text{var}(\ln z) = \sigma^2. \quad (\text{A.4})$$

It follows from the summation result for jointly normal random variables that products and powers of jointly lognormal variables are again lognormal. For example, if  $u$  and  $v$  are lognormal, then  $z = u^\alpha v^\beta$  is also lognormal.

# APPENDIX B

## CALCULUS AND OPTIMIZATION

This appendix reviews the essential elements of calculus and optimization mathematics used in the text.

### B.1 Functions

A function assigns a value that depends on its independent variables. Usually a function is denoted by a single letter, such as  $f$ . If the value of  $f$  depends on a single variable  $x$ , the corresponding function value is denoted by  $f(x)$ . An example is the function  $f(x) = x^2 - 3x$ . We can evaluate this function at  $x = 2$  as  $f(2) = 2^2 - 3 \times 2 = -2$ . Although a function is most properly called by its name, such as  $f$ , it is sometimes convenient, and quite common, to refer to  $f(x)$  as a function, even though  $f(x)$  really is the value of  $f$  at  $x$ .

A function may be defined only for certain numerical values. In many cases, for example, a function is defined only for integer values, in which case the independent variable is usually denoted by  $i, j, k, m$ , or  $n$ . An example is the function  $d(n) = 1/(1+r)^n$ , which is the discount function.

Functions of several variables are also important. For example, a function  $g$  may depend on two variables  $x$  and  $y$ , in which case the value of  $g$  at  $x$  and  $y$  is  $g(x, y)$ . An example is  $g(x, y) = x^2 + 3xy - y^2$ .

Certain types of functions are commonly used in investment science. These include:

1. **Exponential functions** An exponential function is a function of a single variable of the form

$$f(t) = ac^{bt}$$

where  $a, b$ , and  $c$  are constants. Very often the constant  $c$  is  $e = 2.7182818\dots$ , the base of the natural logarithm.

The exponential function also arises when the variable is restricted to be an integer, such as the function  $k(n) = (1+r)^n$ , which shows how capital grows

under compound interest. In this case the function is said to exhibit geometric growth, or to be a geometric growth function.

- 2. Logarithmic functions** The natural logarithm is the function denoted by  $\ln$ , which satisfies the relation

$$e^{\ln(x)} = x.$$

Some important values are  $\ln(1) = 0$ ,  $\ln(e) = 1$ , and  $\ln(0) = -\infty$ .

- 3. Linear functions** A linear function of a single variable  $x$  has the form  $f(x) = ax$ , where  $a$  is a constant. A function  $f$  of several variables  $x_1, x_2, \dots, x_n$  is linear if it has the form

$$f(x_1, x_2, \dots, x_n) = a_1x_1 + a_2x_2 + \dots + a_nx_n$$

for some constants  $a_1, a_2, \dots, a_n$ .

- 4. Inverse functions** A function  $f$  has an inverse function  $g$  if for every  $x$  there holds  $g(f(x)) = x$ . Often the inverse function is denoted by  $f^{-1}$ .

As an example consider the function  $f(x) = x^2$ . This function has the inverse  $f^{-1}(y) = \sqrt{y}$ . Clearly  $f^{-1}(f(x)) = \sqrt{x^2} = x$ . As another example, if  $f$  is the logarithmic function  $f(x) = \ln(x)$ , then the inverse function is  $f^{-1}(y) = e^y$  because  $e^{\ln(x)} = x$ . It is also true that if  $g$  is the inverse of  $f$ , then  $f$  is the inverse of  $g$ . For example, we know that  $\ln(e^x) = x$ .

- 5. Vector notation** When working with several variables it is convenient to regard them as a vector and write, for example,  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ . We then write the value of a function of these variables as  $f(\mathbf{x})$ .

## B.2 Differential Calculus

It is assumed that the reader is familiar with differential calculus. We shall review a certain number of concepts that are used in the text.

- 1. Limits** Differential calculus is based on the notion of a limit of a function. If the function value  $f(x)$  approaches the value  $L$  as  $x$  approaches  $x_0$ , we write

$$L = \lim_{x \rightarrow x_0} f(x).$$

An example is  $\lim_{x \rightarrow \infty} 1/x = 0$ .

- 2. Derivatives** Given a function  $f$ , the derivative of  $f$  at  $x$  is

$$\frac{df(x)}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}.$$

Sometimes we write  $f'(x)$  for the derivative of  $f$  at  $x$ . It is important to know these common derivatives:

- (a) If  $f(x) = x^n$ , then  $f'(x) = nx^{n-1}$ .
- (b) If  $f(x) = e^{ax}$ , then  $f'(x) = ae^{ax}$ .
- (c) If  $f(x) = \ln(x)$ , then  $f'(x) = 1/x$ .

- 3. Higher order derivatives** Higher order derivatives are formed by taking derivatives of derivatives. For example, the second derivative of  $f$  is the derivative of the function  $f'$ . We denote the  $n$ th derivative of  $f$  by  $d^n f/dx^n$ . In the special case of the second derivative we often use the alternative notation  $f''$ .

As an example, consider the function  $f(x) = \ln(x)$ . The first derivative is  $f'(x) = 1/x$ ; the second derivative is  $f''(x) = -1/x^2$ .

- 4. Partial derivatives** A function of several variables can be differentiated partially with respect to each of its arguments. We define

$$\frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_i} = \lim_{\Delta x \rightarrow 0} \frac{f(x_1, x_2, \dots, x_i + \Delta x, x_{i+1}, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{\Delta x}.$$

For example, suppose  $f(x, y) = x^2 + 3xy - y^2$ . Then  $\partial f(x, y)/\partial x = 2x + 3y$  and  $\partial f(x, y)/\partial y = 3x - 2y$ .

We write the total differential of  $f$  as

$$df = \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \dots + \frac{\partial f}{\partial x_n} dx_n.$$

- 5. Approximation** A function  $f$  can be approximated in a region near a given point  $x_0$  by using its derivatives. The following two approximations are especially useful:

$$(a) f(x_0 + \Delta x) = f'(x_0) \Delta x + O(\Delta x)^2$$

$$(b) f(x_0 + \Delta x) = f(x_0) + \frac{1}{2} f''(x_0) (\Delta x)^2 + O(\Delta x)^3$$

where  $O(\Delta x)^2$  and  $O(\Delta x)^3$  denote terms of order  $(\Delta x)^2$  and  $(\Delta x)^3$ , respectively. These approximations apply only to ordinary functions with well-defined derivatives. They do not apply to functions that contain Wiener processes. (See Chapter 13.)

## B.3 Optimization

Optimization is a very useful tool for investment problems. This section reviews only the bare essentials; but these are sufficient for most of the work in the text.

- 1. Necessary conditions** A function  $f$  of a single variable  $x$  is said to have a maximum at a point  $x_0$  if  $f(x_0) \geq f(x)$  for all  $x$ . If the point  $x_0$  is not at a boundary point of an interval over which  $f$  is defined, then if  $x_0$  is a maximum point, it is necessary that the derivative of  $f$  be zero at  $x_0$ ; that is,

$$f'(x_0) = 0.$$

This equation can be used to find the maximum point  $x_0$ .

For example, consider the function  $f(x) = -x^2 + 12x$ . To find the maximum, we set the derivative equal to zero to obtain the equation  $-2x + 12 = 0$ . This has solution  $x = 6$ , which is the maximum point.

A similar result holds when the function  $f$  depends on several variables. At a maximum point (with none of the variables at a boundary point) each of the partial derivatives of  $f$  must be zero. In other words, at the maximum point,

$$\begin{aligned}\frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_1} &= 0 \\ \frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_2} &= 0 \\ &\vdots \\ \frac{\partial f(x_1, x_2, \dots, x_n)}{\partial x_n} &= 0.\end{aligned}$$

This is a system of  $n$  equations for the  $n$  unknowns  $x_1, x_2, \dots, x_n$ .

**2. Lagrange multipliers** Consider the problem of maximizing the function  $f$  of several variables when there is a constraint that the point  $x$  must satisfy the auxiliary condition  $g(x_1, x_2, \dots, x_n) = 0$ . We say that we are looking for a solution to the following maximization problem:

$$\begin{array}{ll}\text{maximize}_x & f(x_1, x_2, \dots, x_n) \\ \text{subject to} & g(x_1, x_2, \dots, x_n) = 0.\end{array}$$

The condition for a maximum can be found by introducing a Lagrange multiplier  $\lambda$ . We form the Lagrangian

$$L = f(x_1, x_2, \dots, x_n) - \lambda g(x_1, x_2, \dots, x_n).$$

We can then treat this Lagrangian function as if it were unconstrained to find the necessary conditions for a maximum. Specifically, we set the partial derivatives of  $L$  with respect to each of the variables equal to zero. This gives a system of  $n$  equations, but there are now  $n + 1$  unknowns, consisting of  $x_1, x_2, \dots, x_n$  and  $\lambda$ . We obtain an additional equation from the original constraint  $g(x_1, x_2, \dots, x_n) = 0$ . Therefore we have a system of  $n + 1$  equations and  $n + 1$  unknowns.

If there are additional constraints, we define additional Lagrange multipliers—one for each constraint. For example, the problem

$$\begin{array}{ll}\text{maximize}_x & f(x_1, x_2, \dots, x_n) \\ \text{subject to} & g(x_1, x_2, \dots, x_n) = 0 \\ & h(x_1, x_2, \dots, x_n) = 0\end{array}$$

can be solved by introducing the two Lagrange multipliers  $\lambda$  and  $\mu$ . The Lagrangian is

$$L = f(x_1, x_2, \dots, x_n) - \lambda g(x_1, x_2, \dots, x_n) - \mu h(x_1, x_2, \dots, x_n).$$

The partial derivatives of this Lagrangian are all set equal to zero, giving  $n$  equations. Two additional equations are obtained from the original constraints. Therefore there are  $n + 2$  equations and  $n + 2$  unknowns.

Some problems have inequality constraints of the form  $g(x_1, x_2, \dots, x_n) \leq 0$ . If it is known that they are satisfied by strict inequality at the solution [with  $g(x_1, x_2, \dots, x_n) < 0$ ], then the constraint is not active and can be dropped from consideration; no Lagrange multiplier is needed. If it is known that the constraint is satisfied with equality at the solution, then a Lagrange multiplier can be introduced, as before. In this case the Lagrange multiplier is nonnegative (that is,  $\lambda \geq 0$ ).

# ANSWERS TO EXERCISES

The answers to all odd-numbered exercises are given here.<sup>1</sup> If the exercise involves a proof, a very brief outline or hint is given.

## Chapter 2

1. (a) \$1,000; (b) \$1,000,000.
3. (a) 3.04%; (b) 19.56%; (c) 19.25%.
5.  $n = 9$  or (equally)  $n = 10$ .
7.  $PV = \$4,682,460$ .
9.  $x < 3.3$ .
11. \$6,948.
13.  $NPV_1 = 29.88$  and  $NPV_2 = 31.84$ ; hence recommend 2.  
 $IRR_1 = 15.2\%$  and  $IRR_2 = 12.4\%$ ; hence recommend 1.
15. (b)  $c = .946$ ,  $r = 5.7\%$ .
17. No inflation applied:  $NPV = -\$435,000$ ; inflation applied:  $NPV = \$89,000$ .

## Chapter 3

1. \$377.12.
3. (a) 95.13 years; (b) \$40,746; (c) \$38,387.
5. No.
7.  $YTM < 9.366\%$ .
9. The annual worths are  $A_A = \$6,449$  and  $A_B = \$7,845$ .

---

<sup>1</sup> Compilation of these answers was the result of a massive project by a number of devoted individuals. We do not guarantee that they are free from errors. Please report errors to the author.

**11.** (a) 850; (b) 1000.

**13.** 91.17.

**15.**  $D = \frac{1+r}{r}, D_M = 1/r.$

**17.**  $dP/d\lambda = -DP.$

**19.**  $C = T^2.$

## Chapter 4

**1.** 7.5%.

**3.**  $P = 65.9.$

**5.** (a)  $f_{t_1, t_2} = [s(t_2)t_2 - s(t_1)t_1]/(t_2 - t_1);$  (c)  $x(t) = x(0)e^{s(t)t}.$

**7.**  $P = 37.64.$

**9.**  $(1+r)^i(1+f_{i,j})^{j-i} = (1+r)^j$  implies  $(1+f_{i,j})^{j-i} = (1+r)^{j-i}$ , which implies  $f_{i,j} = r.$

**11.**  $PV = 9.497.$

**13.** (a) 2.020; (b) 2.02075.

**15.**  $x_1 \approx -13.835,$   $x_2 \approx 30.995.$

**17.**  $a_k = 1/(1+r_{k-1})^2,$   $b_k = 1/(1+r_{k-1}).$

## Chapter 5

**1.** Approximate: projects 1, 2, 5; optimal: projects 1, 2, 3.

**3.**  $NPV = \$610,000$  achieved by projects 4, 5, 6, 7 or 1, 4, 5, 7.

**5.** 16 in lattice, 40 in tree.

**7.** Critical  $d^* = \frac{1}{2}(\sqrt{5} - 1) \approx .618.$  Values  $r = .33$  and  $r = .25$  give  $d = .75$  and  $d = .8,$  so solutions are the same.

**9.** (b)  $PV = \$366,740;$  enhance 2 years, then normal.

**11.** Use hint and solve for  $S.$

## Chapter 6

**1.**  $(2X_0 - X_1)/X_0.$

**3.** (a)  $\alpha = 19/23;$  (b) 13.7%; (c) 11.4%.

**5.** (a)  $(1.5 \times 10^6 + .5u)/(10^6 + .5u);$  (b) 3 million units, 0 variance, 20% return.

**7.** (a)  $\mathbf{w} = (.5, 0, .5);$  (b)  $\mathbf{w} = \left(\frac{1}{3}, \frac{1}{6}, \frac{1}{2}\right);$  (c)  $\mathbf{w} = (0, .5, .5).$

**9.**  $r = \left[\sum_{i=1}^n (1/A_i)\right]^{-1} - 1.$

**11.** (a)  $w_1 = (s_1 - \rho s_2)/F,$  where  $F = (s_1 - s_2)(1 - \rho);$  (b)  $w_1 = 1.$

**13.**  $w_1 = 8/9,$   $w_2 = 1/9.$

## Chapter 7

1. (a)  $\bar{r} = .07 + .5\sigma$ ; (b)  $\sigma = .64$ , borrow \$1,000 and invest \$2,000; (c) \$1,182.
3. (a)  $.1 \leq \bar{r}_M \leq .16$ ; (b)  $.12 \leq \bar{r}_M \leq .16$ .
5.  $\beta_i = x_i \sigma_i^2 \left( \sum_{j=1}^n x_j^2 \sigma_j^2 \right)^{-1}$ .
7. (a)  $A = 1$ ; (b)  $\alpha = \sigma_0^2 / (\sigma_0^2 - \sigma_i^2)$ ; (c) zero-beta point is efficient but below MVP; (d)  $\bar{r}_i = 10\%$ .
9. The identities require simple algebra.
11. (a)  $\beta_a = 0.2$ ,  $\bar{r}_a = 0.04$ ; (b)  $.6364$ ; (c)  $60\%$ .
13. (a)  $r_E = (r_A - wr_B) / (1 - w)$ ; (b)  $\beta_E = \beta_A / (1 - w)$ ; (c) Increase.
15.  $r_f = (\bar{r}_1 + \bar{r}_2) / 2$ .
17. (a)  $w_1 = -1$ ,  $w_2 = 2$ ; (b)  $R = 1.2$ .

## Chapter 8

1. (a)  $11.44\%$ ; (b)  $\sigma = 16.7\%$ .
3. Normalized  $\mathbf{v} = (.217, .263, .360, .153)$ ; eigenvalue =  $311.16$ ; principal component follows market well.
5.  $a = .035693$ ,  $b = .599$ .
7.  $E(r) = 1$ ,  $\sigma = 16.2\%$ .

## Chapter 9

1. (a)  $\sigma(\hat{r}) = \sigma$ ; (b)  $\sigma(\hat{\sigma}^2) = \sqrt{2}\sigma^2 / \sqrt{n-1}$ .
3. *Method:* Index half-monthly points by  $i$ . Let  $r_i$  and  $\rho_i$  be returns for full month and half month starting at  $i$ . Assume the  $\rho_i$ 's are uncorrelated. Then  $r_i = \rho_i + \rho_{i+1}$ . Show that  $\text{cov}(r_i, r_{i+1}) = \frac{1}{2}\sigma^2$ . Find the error in  $\hat{r} = \frac{1}{24} \sum_{i=1}^{24} r_i$ . Ignoring missing half-month terms at the ends of the year, the method gives the same result as the ordinary method.
5. 0.5365.
7. (a)  $13.54\%$ , (b)  $13.91\%$ .

## Chapter 10

1. \$16,317,396.
3.  $\text{VaR}_h = 40 - 100h$ .
5. Use  $\sigma_{1+2} = \sqrt{\sigma_1^2 + 2\sigma_{1,2} + \sigma_2^2} \leq \sqrt{\sigma_1^2 + 2\sigma_1\sigma_2 + \sigma_2^2} = \sigma_1 + \sigma_2$ .
7. All axioms except positive homogeneity are violated.
9.  $\text{CVaR}_h = -50h + 40$ .

- 11.** (a)  $\text{CVaR}_{30\%} = 22/3$ ; (b) In  $\mathcal{P}$ , all columns contain only zeros and a single  $1/3$  and a single  $2/3$ , in any order.

## Chapter 11

- 1.** \$108,610.
- 3.**  $a(x)$ .
- 5.**  $a = (A' - B')/[U(A') - U(B')]$ ,  $b = [B'U(A') - A'U(B')]/[U(A') - U(B')]$ .
- 7.**  $C = (3 + e)^2/16$ ,  $e = 4\sqrt{C} - 3$ .
- 9.**  $b' = b/W$ .
- 11.** \$1,500.
- 13.**  $7/4$ .
- 15.** From hint:  $\bar{R}_i - R = cW[\mathbb{E}(R_M, R_i) - \bar{R}_M R] = cW[\text{cov}(R_M, R_i) + \bar{R}_M(\bar{R}_i - R)]$ . This implies  $\bar{R}_i - R = \gamma \text{ cov}(R_M, R_i)$  for some  $\gamma$ . Apply to  $R_M$  to solve for  $\gamma$ .
- 17.**  $P = \mathbb{E}\left(\frac{d}{R^*}\right) = \mathbb{E}\left(\frac{Rd}{RR^*}\right) = \frac{1}{R}\mathbb{E}\left(\frac{Rd}{R^*}\right) = \frac{\hat{\mathbb{E}}(d)}{R}$ .

## Chapter 12

- 1.** \$442.02.
- 3.** 5%.
- 5.** There is no cash flow at  $t = 0$ . At  $T$  the flow is  $S/d(0, M) + \sum_{k=0}^{M-1} c(k)/d(k, M) - F$ , which must be zero.
- 7.** -\$100.34.
- 9.**  $\frac{S(0)}{r}[1 - e^{-rT}]$ .
- 11.** (a)  $V_{i-1}(r_i) = 1 - d(i-1, i)$ ; (b)  $V_0(r_i) = d(0, i-1) - d(0, i)$ ; (c)  $1 - d(0, M)$ .
- 13.** (a) \$3.971 million; (b) 8.64%.
- 15.** -131,250 lb orange juice;  $\sigma_{\text{new}} = .714\sigma_{\text{old}}$ .
- 17.** Short \$163,200 Treasury futures.
- 19.** Proof based on  $\text{cov}(x, y^2) = \mathbb{E}(xy^2) - \mathbb{E}(xy)\mathbb{E}(y) = 0$ . Both  $\mathbb{E}(xy^2)$  and  $\mathbb{E}(y)$  are zero by symmetry.

## Chapter 13

- 1.** Assuming  $\Delta t$  small,  $p = .65$ ,  $u = 1.106$ ,  $d = .905$ ; without small  $\Delta t$  approximation,  $p = .64367$ ,  $u = 1.11005$ ,  $d = .90086$ . Probabilities of nodes (from the top with small approximation) are .179, .384, .311, .111, .015.
- 3.** (a) Use  $(v_1 - v_2)^2 \geq 0$ ; (b) 15% and 9.54%; (c) arithmetic for simple interest, geometric for compound. Usually geometric is best.
- 5.**  $\text{var}(u) = e^{2\bar{w}+\sigma^2}(e^{\sigma^2} - 1)$ .

7.  $dG = (\frac{1}{2}a - \frac{1}{8}b^2)G dt + \frac{1}{2}bG dz.$   
 9.  $dF = (\mu - r)F dt + \sigma F dz.$   
 11. (a)  $S(t)e^{\mu(T-t)}$ ; (b)  $\sigma W(t)dz$ ; (c)  $-\sigma^2 t$ ; (d)  $e^{\sigma^2 t} - 1$  and  $\sigma^2 t$ .  
 13. To first order both have expected value  $S(t_{k+1}) = (1 + \mu \Delta t)S(t_k)$ .

## Chapter 14

1. Cost is nonnegative.  
 3. \$8.4 million.  
 5. \$2.83 American, \$2.51 European.  
 7.  $C(S, T) \geq \max[0, S - KB(T)] \rightarrow S$  as  $T \rightarrow \infty$ . Clearly  $C(S, T) \leq S$ . Hence in the limit  $C = S$ .  
 9. \$7.  
 11. Offer is close: low by about .3%.  
 13. Almost identical! One-month interval: \$4.801; half-month: \$4.796.  
 15. \$6.58.  
 17. (a)  $r_f = 0$ ; (b)  $q_h = 6/10$ ,  $q_t = 3/10$ ,  $q_e = 1/10$ .  
 19. (a) Butterfly (triangle); (b) 10%; (c) \$6.50.

## Chapter 15

1. \$2.57.  
 3.  $\sigma = .251$ .  
 5.  $C(63) = \$6.557$ ,  $\Delta = .759$ ,  $\Theta = 6.02$ .  
 7.  $\Gamma = \frac{\partial \Delta}{\partial S} = \frac{\partial N(d_1)}{\partial S} = N'(d_1) \frac{\partial d_1}{\partial S} = \frac{N'(d_1)}{S\sigma\sqrt{T}}$ . For  $\Theta$  use  $\Gamma$  and Exercise 6.  
 9. (a)  $dV = (V_t + V_x \eta(\theta - x) + \frac{1}{2}V_{xx}\sigma^2) dt + V_x \sigma dz$ ;  
     (b)  $rV(x, t) = V_t + V_x rx + \frac{1}{2}V_{xx}\sigma^2$ .  
 11.  $P_r = KTe^{-rT}[N(d_2) - 1] = -KTe^{-rT}N(-d_2)$ .  
 13.  $a = -\text{cov}(x, y)/\text{var}(y)$ .  
 15. (a) \$.53; (b) \$2.04.  
 17. \$42.42 million.

## Chapter 16

1. (a) 91.72; (b) 90.95.  
 3. Do backward evaluation on futures price lattice.  
 5. 6.00, 6.15, 6.29, 6.44, 6.59, 6.74, 6.89, 7.05, 7.19, 7.35 percent.

7. 7.67, 8.829, 9.799, 10.66, 11.3, 11.93 are  $a_0$  through  $a_5$ .  
 9. \$162,800.  
 11.  $F(t) = r - \frac{1}{2}at + \frac{1}{6}\sigma^2t^2$ .  
 13. \$930.07.

## Chapter 17

1. 28.5%.  
 3. \$670,230.  
 5.  $1/\lambda$ .  
 7.  $[1 - e^{-(r+\lambda)T}]F$ .  
 9.  $P[\tau \geq t] = P[q^{-1}(U) \leq t] = P[U \leq q(t)]$ .  
 11. (a)  $p(1-p)$ ; (b)  $\sigma_{ab} = \sqrt{p(1-p)/n}$ ; (c) 1 million.  
 13.  $P_A = 104.160, P_B = 72.769$ .

## Chapter 18

1.  $\gamma = \frac{1}{4}$ .  
 3.  $\max\left\{\frac{1}{2}\ln[2\alpha + (1-\alpha)] + \frac{1}{2}\ln[\alpha/2 + (1-\alpha)]\right\}$  gives  $\alpha = \frac{1}{2}$ .  
 5. (a)  $\alpha_k = p_k - p_n r_n / r_k$  for  $k < n$ ; (b)  $\alpha_1 = \frac{5}{18}, \alpha_2 = 0, \alpha_3 = \frac{1}{18}$ .  
 7. Dow Jones average outperforms Mr. Jones.  
 9.  $w_1 = .259, w_2 = .185, w_3 = .556$ .  
 11. (a)  $\alpha = (3p - R)/(3 - R)$ ; (b)  $q = R/3$ ; (c)  $\alpha = 0$ .  
 13. (b)  $\mu_i - r = (1 - \gamma)\text{cov}(P_{\text{opt}}, S_i)$ .

## Chapter 19

1. (a) Yes, use portfolio weights  $\frac{1}{3}, \frac{2}{3}$  to get 1.2 risk free; (b) yes, use weights  $-\frac{1}{2}, \frac{1}{2}$ .  
 3. (a) and (b) \$.8678.  
 5.  $q_{11} = .1, q_{12} = .36, q_{21} = .4, q_{22} = .14$ .  
 7.  $S = \$16.81, \sigma = 20.6\%$ .  
 9. Car B preferred by certainty equivalent difference of \$370.74.  
 11.  $V_1 = \frac{E(x_{2|1})}{1 + r + \beta_2(\bar{r}_2 - r)}, \beta_2 = \frac{\text{cov}(x_{2|1}/V_1, r_2)}{\sigma_{r_2}^2}$

# INDEX

- Accrued interest, 52–53, 70
- “Actual” curve, 243–248, 250
- Actual dollars, 35
- Additive model, 353–355
- Adjustable-rate mortgages, 10, 43, 46, 457–461
- Affine processes, 476–477
- After-tax rate, 33–34
- Aggregate drift, 464
- American option, 375–376
- Amortization, 49–51
- Annual percentage rate (APR), 49–50
- Annual worth, 51–52
- Annuities, 46–49
- Annuity formulas, 49
- APT (arbitrage pricing theory), 223–227, 231
- Arbitrage
  - description of, 4–5, 550
  - possibilities for, 293, 454
  - type A, 291–292
  - type B, 292–293
- Arbitrage argument, 82–83, 454
- Arbitrage bounds, 296–297
- Aristotle, 374–375
- Arrow-Pratt coefficient, 284
- “As you like it” option, 432
- Asian option, 433
- Asset return, 144–147
- Assets, 325–326, 549
  - at the money, 378
- Atom, 259, 272
- Automobile purchases, 31–32
- Average value at risk (VaR), 270
- Axioms
  - Black-Scholes, 438–440
  - Pricing, 438–440, 541, 572–575
  - Risk-measure, 269
  - Utility function, 281
- Backward equation, 475–476
- Backwardation, 336
- Balloon payment, 46
- Banker’s acceptance, 44
- Bell-shaped curve, 261
- Benefit-cost ratio, 109
- Bermudan option, 432
- Beta, 186–187
- Betting wheel, 154–155
- Binomial form, 435–436
- Binomial lattice
  - in dynamic choice, 117–119
  - interest rate model, 451–454
  - option model, 389–393, 429–431
  - stock model, 351–353, 368–369, 391
- Binomial options theory, 383–386
- Binomial tree, 117–119
- Black and Karasinski model, 474
- Black-Derman-Toy model, 465, 473–474
- Black-Scholes equation, 410–414, 540–541, 573–574
  - derivation, 410–414
  - and log-optimal, 540–541

- generalized, 573–575
- by axioms, 438–440
- Blur**
  - of  $\alpha$ , 241–242
  - of history, 236
  - of mean, 238–240
- Bond price formula**, 55
- Bonds**
  - accrued interest, 52–53, 70
  - callable, 45, 449
  - CMOs, 469
  - corporate, 45
  - coupon payments, 44–45, 52–59, 61, 95–96
  - derivatives, 455, 506
  - description of, 52–53
  - embedded options, 449
  - face value, 52
  - floating-rate, 95–96
  - futures, 449
  - general obligation, 45
  - junk, 54, 501–503, 507–508
  - long, 59
  - municipal, 45
  - par, 52, 57
  - price formula, 55
  - price sensitivity, 63, 97
  - price-yield curve, 55–58
  - putable, 449
  - quality of, 45, 76
  - revenue, 45
  - short, 59
  - U.S. Treasury, 44, 52
  - yields, 76–77
  - zero-coupon, 45
- Brownian motion**, 360, 363, 365, 436–437
- Bull spread**, 404
- Butterfly spread**, 380
- Buying price**, 566–572
- Calculus**, 439
- Call option**
  - and Black-Scholes formula, 423–425
  - concept of, 374–376
  - formula for, 414–416
  - perpetual, 405, 412, 433
  - pricing of, 382–386
  - value of, 377–380
- Callable bond**, 45, 449
- Capital asset pricing model (CAPM)**
  - and APT, 227
  - development of, 180
  - as factor model, 220–223
  - investment implications, 190–191
  - performance evaluation, 191–194
  - as pricing formula, 194–197
  - projection pricing, 200–202
  - security market line, 187–189
  - systematic risk, 189–190
  - theorem of, 184–187
- Capital budgeting**, 108–110, 111–113
- Capital market line**, 182–184
- Capitalization weights**, 181
- CAPs**, 432
- Carrying charges**, 323, 345–346
- Cash flow**
  - description of, 1–3
  - free, 132–134
  - in graphs, 119–120
- Cash flow stream**
  - description of, 1–3
  - deterministic, 9
  - general, 11
  - single-period random, 10
- Cash matching**, 114–116
- Cauchy–Schwarz inequality**, 204
- CDS (credit default swap)**, 8, 483–484, 490, 506–508
- Certain value**, 149
- Certainty equivalent**, 196–197, 284–285, 567–571
- Certificate of Deposit (CD)**, 43
- Characteristic line (or equation)**, 221–223
- “Chooser” options, 432
- CIR (Cox, Ingersoll, and Ross) model**, 474, 497–498
- Coherent characterization theorem**, 272–273
- Coherent risk measures**, 269–270

- Collateralized debt obligations (CDO's), 509–511
- Collateralized mortgage obligations (CMOs), 469–473
- Commercial paper, 44
- Committee on Uniform Securities Identification Procedures (CUSIP), 45
- Commodity swap, 327–329
- Common stocks, beta of, 187
- Comparables, method of, 203–204
- Comparison principle, 4
- Complete market, 302
- Complexico gold mine, 125–127, 138, 395–396, 406, 458, 577
- Compound interest, 16–19, 97–98
- Compound options, 432
- Compounding, 17–19, 23, 97–98, 248
- Concave utility, 282–283
- Condition number, 242
- Conditional tail expectation (CTE), 270
- Conditional value at risk (CVaR), 270–272
- Confidence level, 258
- Consols, 47
- Constant dollars, 35
- Constant-growth dividend, 131
- Contango, 336
- Continuous compounding, 18–19, 23
- Continuous-time
  - growth, 528–531
  - pricing axioms, 572–575
- Control variate, 427, 445
- Convenience yield, 325
- Convergence, 332, 435
- Convexity, 68–69, 274–275
- Copulas, 505
- Corporate bonds, 45
- Correlated defaults, 503–505
- Correlation coefficient, 152
- Correlation pricing, 203–206
- Cost of capital, 29
- Cost of carry, 322–324
- Cost of risk, 440
- Counterparty risk, 8, 483
- Coupon payments, 44–45, 52–59, 61, 95–96
- Coupon rates, 56, 62, 70
- Coupons, 491–492
- Covariance, 580
  - definition of, 150–152
  - in portfolio return, 156–159, 169–172
- Covariance matrix, 581
- Cox, Ingersoll, and Ross (CIR) model, 474, 497–498
- Cox processes, 495, 497–498
- Credit default swap (CDS), 8, 483–484, 490, 506–508
- Credit derivatives, 505–511
- Credit rating methods, 504
- Credit risk, 258, 498–500
- Credit spread, 486–487
- Cross-ratio options, 432
- CTE (conditional tail expectation), 270
- Cumulative probability distribution, 259, 414
- Current yield (CY), 59–60
- Curse of dimensionality, 120
- CUSIP (Committee on Uniform Securities Identification Procedures), 45
- CVaR (conditional value at risk), 270–272
- CY (current yield), 59–60
- Cycle problems, 31–33
- Debt subordination, 45
- Default, 43, 486, 490, 492, 503–505
- Delta, 417–419, 443
- Delta property, 567–568
- Demand deposit, 43
- Depreciation, 34
- Derivatives
  - asset, 10
  - bond, 455
  - credit, 505–511
  - and differential calculus, 584–585
  - higher order, 585
  - interest rate, 448–450

- partial, 585
- and utility functions, 284
- Deterministic cash flow stream, 9
- Differential calculus, 584–585
- Digital option, 432
- Dimensionality, curse of, 120
- Direct simulation, 498–499
- Discontinuous options, 432
- Discontinuous value, 268
- Discount factor, 20–21, 79, 89–90
- Discounted growth formula, 131
- Discrete-time compounding, 97–98
- Diversifiable risk, 217
- Diversification, 157–159
- Diversification failure, 266–267
- Dividends
  - discount model, 130–131
  - and options, 391
  - and storage costs, 435–436
- Dollars, 35
- Double lattice, 555–560
- Doubly stochastic process, 495
- Dow Jones Average, 543–544
- “Down and out” options, 432–433, 488
- Drift, 187, 437, 464, 537
- Duration
  - and sensitivity, 62–64
  - definition, 59
  - Fisher-Weil, 96–97
  - formula for, 59
  - Macaulay, 60–61
  - modified, 63, 65, 68–70
  - portfolio, 64–65
  - quasi-modified, 98
- Dynamic cash flow processes, 116–120
- Dynamic hedging strategy, 418
- Dynamic model, 117
- Dynamic programming, running, 120–127
- Dynamics
  - expectations, 88–92
  - of financial markets, 5
  - interest rate, 469–474
- Early default, 490
- Early exercise, 382–383, 388–389
- Effective interest rate, 17–18
- Efficient frontier, 164, 531–532
- Elementary prices, 462
- Elementary state securities, 302
- Embedded bond options, 449
- Equal and opposite hedge, 338–339, 344–345, 348
- Equilibrium, 180, 182, 250–252
- Equivalent streams, 23
- Equivalent utility functions, 281–282
- Errors, 214
- Estimation errors, 242–248
- Eurodollars, 44
- European option, 375–376
- Excess returns, 186, 221–222
- Exchange options, 432
- Exercise, early, 382–383, 388–389
- Exotic options, 431–433
- Expectations dynamics, 88–92
- Expectations hypothesis, 85–86
- Expected shortfall, 270
- Expected spot price, 335–336
- Expected value, 148, 580
- Exponential functions, 583–584
- Exponential growth curve, 19
- Exponential utility, 280–281, 567–568
- External factors, 220
- Extracted factors, 220
- Face value, 52
- Factor loading, 215
- Factor model
  - approach, 213–214
  - CAPM as, 220–223
  - factor selection, 219–220
  - multifactor, 219
  - portfolio parameters, 215–218
  - projection pricing with, 227–229
  - single-factor, 214–215
- Factor price, 225
- Fat tail, 264
- FCF (free cash flow), 132–134
- Feasible region, 161–165, 531–536
- Film venture, 295–296
- Financial instrument, 42
- Finite state model, 301–304

- Finite-difference method, 427–429, 435  
 Firm value, 130, 220, 484  
 First passage times, 487–488  
 Fisher-Weil duration, 96–97  
 Fishing problem, 123–125, 138  
 Fixed-income securities, 42–43  
 Fixed-proportion, 519  
 Floating-rate bonds, 95–96  
 Forward contract, 315, 318–319,  
     326–327, 455–457, 508  
 Forward equation, 461–464  
 Forward interest rates, 82–85, 319  
 Forward market, 319  
 Forward price, 318, 319–322  
 Forward start options, 432  
 Forward value, 426, 438  
 Free cash flow (FCF), 132–134  
 Function, 583–584  
 Future-forward equivalence, 332–335  
 Future value, 21–22  
 Futures contract, 329–331, 457  
 Futures market, 330  
 Futures options, 391–393  
 Futures prices, 332–335
- Gamma, 418, 444, 545  
 Gaussian random variables, 580–581  
 General cash flow streams, 11  
 General obligation bonds, 45  
 Geometric Brownian motion, 363, 365  
 Geometric growth, 16–17  
 Geometric mean, 192, 371  
 Gold mine  
     Complexico, 125–127, 138,  
         395–396, 406, 458, 577  
     Simplico, 30–31, 81, 393–395,  
         397–399, 402, 559–560  
 Gordon formula, 131  
 Graph, cash flow in, 117, 119–120  
 Growth efficiency proposition,  
     527–528
- HARA utility, 308  
 Harmony theorem, 128–130, 139,  
     198–200
- Heath, Jarrow, and Morton model, 478  
 Hedge  
     description of, 7–8  
     minimum-variance, 337–340  
     nonlinear, 341–344  
     optimal, 340–341  
     perfect, 336  
 Higher order derivatives, 585  
 Ho-Lee model, 464–465, 473,  
     475–476, 479  
 Homogenous process, 495  
 Hull and White model, 474
- Ideal bank, 21, 23  
 Idiosyncratic risk, 189–190  
 Indifference price, 298, 578  
 In the money, 378  
 Immunization, 65–68, 98–100,  
     467–469  
 Implied forward rates, 83  
 Implied volatility, 422–424, 465–467  
 Indenture, 45  
 Independent variable, 150, 580  
 Index fund, 191  
 Inflation, 34–36  
 Inhomogeneous process, 495  
 Inner product, 201  
 Insurance  
     hedging as, 7  
     portfolio, 422  
 Intensity method, 493–494  
 Intercept, 215  
 Interest  
     compound, 16–19, 97–98  
     effective, 17–18  
     nominal, 17–18, 23  
     real, 35  
     simple, 15–17  
 Interest duration, 60. *See also* duration  
 Interest rate  
     derivatives, 448–450  
     dynamics, 469–474  
     forward, 319  
     swap, 329  
 Intermediate receipts, 496–497

- Internal rate of return (IRR)  
 alternative ranking, 28–30  
 main theorem, 24–26  
 project evaluation, 548
- Invariance theorem, 91–92
- Inverse floater, 481
- Inverse functions, 584
- Inverted yield curve, 77
- Investment assets, 325–326
- Investment grade, 54
- Investment wheel, 517–519
- Ito process, 362, 363–364
- Ito's lemma, 364, 366–368
- James-Stein Shrinkage estimator, 249
- Jensen's index, 193
- Joint probability, 579
- Junk bonds, 54, 501–503, 507–508
- Kelly rule of betting, 521–522
- Knockout options, 432
- Lagrange multipliers, 165–166, 586–587
- Lagrangian, 165, 294, 532
- Lattice methods. *See also* binomial lattice; trinomial lattice  
 double, 555–560  
 for credit derivatives, 500–503  
 for defaultable bonds, 488–489
- Law of large numbers, 520
- LEAPS (Long-term Equity AnticiPation Securities), 432
- Leveling, 458
- Leverage, 242–243
- LIBOR, 449n1, 478, 506
- Limits, 584
- Linear functions, 584
- Linear independence, 200
- Linear pricing, 196, 291, 317–318, 399–401
- Linearity, 149
- Liquidity preference, 86–87
- Loans, 506
- Log utility, 519–525
- Logarithmic function, 281, 584
- Logarithmic performance, 520–521
- Lognormal prices, 355–356, 363
- Lognormal variables, 358–359, 581
- Log-optimal pricing, 299–301, 536–539, 540–541
- portfolio, 533
- pricing, 299–301, 554
- pricing formula (LOPF), 536–539
- Black-Scholes, 540–541, 573–574
- Log-optimal strategy, 525–526
- Long, 318
- Long bond, 59
- Lookback option, 432
- Loss tolerance, 258
- Macaulay duration, 60–61
- Machine replacement, 32–33
- Margin, 376
- Margin account, 330
- Margin call, 331
- Marginal price, 298
- Market capitalization weights, 181
- Market data, 539–540
- Market equilibrium, 180–182
- Market forward rates, 83
- Market portfolio, 181
- Market price of risk, 440, 574
- Market risk, 258, 265–266
- Market segmentation, 87
- Market space, 200
- Market uncertainty, 560–562
- Marking to market, 330
- Markowitz model, 164–167, 191, 230
- Martingale pricing, 437–438
- Master asset, 146
- Maximum tangent, 245–248
- Mean blur, 238–239, 241
- Mean return, 156
- Mean reversion, 474
- Mean value, 148
- Mean-standard deviation, 155–156
- Mean-variance analysis, 6, 10, 143
- Merton model, 484–487

- Minimum norm pricing, 202–203
- Minimum-variance hedge, 337–340
- Minimum-variance point (MVP), 163
- Minimum-variance set, 162–163
- Modified duration, 63, 65, 68–70
- Money market, 19–20
- Money market instruments, 44
- Monotonicity, 269
- Monte Carlo simulation, 426–427, 434–435, 445, 498
- Mortgage-backed securities, 46, 449
- Mortgages
  - adjustable-rate, 457–461
  - biweekly, 72
  - features of, 449
  - variations of, 46
- Multifactor models, 217, 220, 474, 477
- Multiperiod fallacy, 229–230
- Multiperiod options, 386–389
- Multiperiod securities, 548–550
- Multiplicative model, 355–356
- Municipal bonds, 45
- Mutual fund, 169
  
- Negatively correlated, 151
- Net present value (NPV), 27–28
- Netting, 327
- No arbitrage opportunities, 293, 454
- Nodes, 117
- Nominal dollars, 35
- Nominal interest, 17–18, 23
- Nondiversifiable risk, 217
- Nonlinear risk, 341–344
- Nonmarket risk, 258
- Nonnegativity, 149, 167–168
- Nonsatiation, 164
- Nonsystematic risk, 189–190
- Norm, 201
- Normal backwardation, 336
- Normal price distribution, 354–355
- Normal probability
  - distribution-density, 261, 414
- Normal random variables, 290–291, 354–356, 486, 580–581
  
- Normal returns, 290–291
- Notional principal, 327
- NPV (net present value), 27–28
  
- Oil well (Rapido), 563, 566, 571
- Oil venture, 184, 195
- One-fund theorem, 173
- Operational calculus, 439–440
- Optimal hedging, 340–341
- Optimal portfolio, 113–116, 303, 533–538
- Optimal pricing, 552–555
- Optimization, 108–109, 113, 585–587
- Option
  - definition, 374
  - strike price, 375
  - exercise price, 375
  - early exercise, 382–383, 388–389
  - synthetic, 419–421
  - American, 375–376
  - binomial lattice for, 383–386, 389–393, 429–431
  - bond, 449
  - call, 374–380, 382–386, 405, 412, 414–416, 423–425, 433
  - combination, 380–381
  - concept of, 375–377
  - and dividends, 391
  - “down and out,” 432–433, 488
  - European, 375–376
  - exercise of, 374
  - exotic, 431–433
  - formula, 575–576
  - futures, 391–393
  - multiperiod, 386–389
  - pricing of, 383–386
  - putable, 389–390
  - real, 397–399, 573
  - time value of, 379
  - value of, 377–380
- Orthogonality, 201–204
- Out of the money, 378
  
- Par, 52, 57
- Parameter

- estimation, 235
- estimation errors, 242–248
- mean blur, 238–239, 241
- period-length effects, 236
- shrinkage estimation, 249
- Parameter estimation, 357, 358
- Partial derivatives, 585
- Path dependency, 433, 458
- Pay-later option, 445
- Perfect hedge, 336
- Perfect market, 317
- Performance evaluation, 191–194
- Period-length effects, 236–238
- Perpetual annuity, 47–48
- Plain vanilla cap, 450, 467
- Plain vanilla swap, 327
- Points, 49
- Poisson process, 493–494
- Portfolio
  - beta of, 187
  - choice theorem, 293–296
  - curve, 185
  - delta, 418
  - diagram of, 159–161
  - dynamic, 529–531
  - effective, 166
  - insurance, 422, 441
  - log-optimal, 533
  - market, 181
  - mean, 156
  - optimal, 113–116
  - pricing, 291
  - replicating, 385, 419–421
  - return, 146–147, 156–157
  - selection problem, 8
  - strategies, 549–550
  - to mean add variance
    - (mean-variance), 156
  - well-diversified, 225–226
- Portfolio optimization, 113
- Position, 258, 318
- Positive homogeneity, 269
- Positive state prices, 302–306
- Positively correlated, 151
- Power utility function, 281
- Premium, 7, 374
- Present value, 20–24, 92–95, 547–548
- Present worth, 27
- Pricing axioms, 438–440, 541, 572–575
- Pricing form of CAPM, 195
- Price of risk, 225, 574
- Price sensitivity formula, 63–64
- Prices, forward, 319–322
- Price-yield curve, 55–58
- Pricing. *See also* capital asset pricing model (CAPM)
  - axioms, 438–440, 541, 572–575
  - buying analysis, 566–572
  - correlation, 203–206
  - as investment problem, 6–7
  - indifference price, 298, 578
  - linear, 291–293, 317–318, 399–401
  - log-optimal, 299–301, 536–539, 540–541
  - martingale, 437–438
  - minimum norm, 202–203
  - optimal, 552–555
  - projection, 202
  - risk-neutral, 304–306, 401–402
  - zero-level, 297–299, 561
- Pricing principles, 316–318, 402–403
- Principal, 15
- Private uncertainty, 560–562
- Probability. *See also* cumulative probability distribution; default probability; risk-neutral probability
  - joint, 579
  - density, 148
  - distribution, 579
  - structure, 579
  - survival, 494–495
  - theory, 579–582
- Project choice, 198–200
- Projection pricing, 200–202, 227–229, 574
- Projection theorem, 202
- Projection pricing with factors, 200–207, 227–229
- Pumping, 522–523
- Pure investment, 8–9

- Putable bonds, 449
- Put-call parity, 380–381
- Put option, 374, 389–390
  
- Quadratic program, 168
- Quadratic utility, 281, 288–290
- Quality ratings, *See also* Rating methods
- Quantile, 259
- Quasi-modified duration, 98
- Questionnaire method, 288
  
- Random returns, 152–155
- Random variables, 147–152
- Random walk, 359–360
- Rapido oil well, 563, 566, 571
- Rate of return, 144
- Rating methods, 53–54, 492–493, 504
- Ratio theorem, 557–558
- Real dollars, 35
- Real interest rate, 35
- Real investments, 393–396
- Real options, 11, 397–399, 573
- Real stock distributions, 356
- Rebalance, 68, 418
- Reduced-form approach, 493
- Reduced-form methods, 505
- Rendleman and Bartter model, 473
- Replicating portfolio, 385, 419–421
- Repo rate, 322
- Returns
  - asset, 144–147
  - excess, 186, 221–222
  - normal, 290–291
  - portfolio, 146–147, 156–157
  - random, 152–155
  - rate of, 144
  - total, 144
  - uncertain, 10
- Revenue bonds, 45
- Risk, market price of, 440
- Risk assessment, 8, 267–268
- Risk aversion, 5–6, 164, 282–285
- Risk management, 8
- Risk neutral, 280
  
- Risk quiz, 289
- Risk-free asset, 165, 171–175, 532–536, 550
- Risk measure, 269
- Risk-neutral pricing, 304–306, 401–402, 477, 550–551, 557
- Risk-neutral probability, 385, 424–425, 551
- Risk-neutral utility, 283, 305
- Risk-neutral valuation, 416–417
- Running dynamic programming, 120–127
- Running amortization, 50–51
- Running present value, 92–95
  
- Savings deposits, 43
- Security, 42, 301
- Security market line, 187–188
- Self-financing, 414
- Sensitivity, and duration, 62–64
- Sequential CMOs, 469
- Seven-ten rule, 16
- Sharpe ratio, 194
- Short, 318
- Short bond, 59
- Short rate, 90–91
- Short rate lattice, 451–452
- Short sales, 144–146
- Short-term credit spread, 494
- Short-term risk free rates, 550
- Shrinkage estimators, 249
- Simple interest, 15–17
- Simplico gold mine, 30–31, 81, 393–395, 397–399, 402, 559–560
- Simulation, 365–366, 426–427, 498–500
  - credit risk, 498–500
  - Monte Carlo, 426–427, 434–435, 445, 498
  - of wheel, 525
  - price process, 365–366
- Single-factor models, 474, 477
- Sinking funds, 45
- Specific risk, 189–190
- Spot market, 319

- Spot price, 335–336
- Spot rates, 78–79, 81–82, 88–89
- Standard deviation, 149
- State prices
  - elementary, 302
  - positive, 302–306
- States, 117, 301
- Stationary process, 519
- Stein’s paradox, 249
- Stochastic intensity model, 495–496
- Stock price process, 362–366
  - additive, 355–356
  - binomial 351–353, 368–369
  - Ito form, 363–365
  - log normal, 363
  - multiplicative, 355–356
- Storage costs, 435–436
- Strike price, 375
- Structural methods, 484, 504–505
- Subadditivity, 267, 269
- Survival probability, 494–495
- Swaps, 327–329, 450
- Swaptions, 450
- Synthetic derivative, 419–421
- Systematic risk, 189–190, 217
  
- Tail value at risk (TVaR), 270
- Tails, 356
- Taxes, 33–34
- Term structure
  - explanations, 85–87, 391
  - Ho-Lee model, 475–476, 479
  - implied, 452–454
  - matching, 464–467
  - of volatility, 425
  - theory, 78–82
- Theorem
  - Black-Scholes equation, 412–413
  - CAPM, 184
  - coherent characterization, 273
  - correlation pricing, 204–205
  - floating-rate value, 95–96
  - forward contract, 326–327
  - forward price, 320–321, 323
  - futures-forward equivalence, 332–334
  - harmony, 128–130, 198–200
  - invariance, 91–92
  - IRR, 25–26
  - Ito’s lemma, 367
  - LOPF, 537
  - minimum-variance hedging, 337–338
  - no arbitrage opportunities, 297
  - normal distribution, 261–262
  - one-fund, 173
  - portfolio choice, 293–294
  - portfolio diagram lemma, 160–161
  - positive state prices, 303
  - present value, 24
  - projection price relation, 202
  - ratio, 557–558
  - risk-neutral probabilities, 551
  - simple APT, 225
  - two-fund, 168–169, 532
  - “Theory” curve, 243–248, 250
  - Theta, 419, 444
  - “Think” curve, 243–248, 250
  - Three-principles, 402–403
  - Three views, 243
  - Tight markets, 324–325
  - Tilting, 250
  - Time value of money, 15, 36
  - Time value, of options, 379
  - Total return, 144
  - Tracking, 178
  - Trading strategy, 549, 550
  - Tranches (of CMO bonds), 469
  - Translation invariance, 269
  - Treasury bills, 52
  - Tree cutting example, 27–29, 38, 561–562
  - Trinomial lattice, 119, 429–431
  - TRORS’s (total rate-of-return swaps), 508
  - TRS’s (total return swaps), 508
  - TVaR (tail value at risk), 270
  - Two-fund theorem, 168–169, 532

- Type A arbitrage, 291–292
- Type B arbitrage, 292–293
- Uncertain returns, 10
- Uncorrelated, 151
- Underlying security, 316
- Universal, 307, 309, 574
- U.S. Government Securities, 44–45
- Utility, measurement of, 285–286
- Utility functions
  - and derivatives, 284
  - axioms, 281
  - concave, 282
  - exponential, 280–281
  - equivalent, 281–282
  - increasing, 280
  - logarithmic, 281
  - mean-variance criterion, 288–290
  - power, 281
  - quadratic, 281, 288–290
  - specification of, 285–288
- Valuation of firm, 130–132
- Value at risk (VaR)
  - capital requirement, 260–261
  - computation of, 261–266
  - conditional, 270–272
  - criticisms of, 266–268
  - definition, 258–260
  - properties of, 260
- Vanilla swap, 327
- Variance, 149–150, 580
- Variance of a sum, 152
- Variance of portfolio return, 156–157
- Variance reduction, 427
- Vasicek model, 474
- Vector notation, 584
- Volatility parameter, 464
- Volatility pumping, 522, 528, 543
- Volatility smiles, 422–425
- Volatility surface, 425
- Weight
  - capitalization, 181, 230
  - in total investment, 146
- Well-diversified portfolios, 225–226
- Wheel
  - betting, 154–155
  - of fortune, 153
  - investment, 517–519
- When to cut a tree, 27–29, 561–562
- White noise, 362
- Wiener process, 360–362
  - generalized, 361
- Write an option, 376
- Yield, 54–55
- Yield curve, 55–58, 76–77
- Yield to call (YTC), 59
- Yield to maturity (YTM), 54
- Yield-based options, 432
- Zero-beta assets, 208
- Zero-coupon bond, 45, 64, 66, 79, 81, 91–92, 461–466
- Zero-level pricing, 297–299, 307, 561
- Zero-one programming problem, 109
- Zero-one variable, 108–109, 111–112



Printed in the USA/Agawam, MA  
June 26, 2019



705894.011











*Investment Science*, Second Edition, provides thorough and highly accessible mathematical coverage of the fundamental topics of intermediate investments, including fixed-income securities, capital asset pricing theory, derivatives, and innovations in optimal portfolio growth and valuation of multi-period risky investments.

Eminent scholar and teacher David G. Luenberger, known for his ability to make complex ideas simple, presents essential ideas of investments and their applications, offering students the most comprehensive treatment of the subject available.

#### NEW TO THIS EDITION

- **Three new chapters:** Risk Measures, Credit Risk, and Data and Statistics
- **Updated content and expanded coverage of many topics,** including the capital asset pricing model, projection pricing, the Black-Scholes equation, computational methods, real options, the characterization of volatility, parameter estimation, and portfolio design
- **New exercises** reflecting advances in theory and practice provide opportunities to explore a wide range of concepts

#### ABOUT THE AUTHOR

**David G. Luenberger** is Professor of Management Science and Engineering, Stanford University.

**OXFORD**  
UNIVERSITY PRESS  
[www.oup.com/us/he](http://www.oup.com/us/he)

Cover Design: M. Laseau  
Cover Image: iStockphoto

