

*Dynamic Programming and Optimal Control*

*Volume II*

*Approximate Dynamic Programming*

FOURTH EDITION

Dimitri P. Bertsekas

Massachusetts Institute of Technology

WWW site for book information and orders

<http://www.athenasc.com>



Athena Scientific, Belmont, Massachusetts

**Athena Scientific**  
**Post Office Box 805**  
**Nashua, NH 03061-0805**  
**U.S.A.**

**Email: [info@athenasc.com](mailto:info@athenasc.com)**  
**WWW: <http://www.athenasc.com>**

Cover Design: Ann Gallager, [www.gallagerdesign.com](http://www.gallagerdesign.com)  
Cover photography: Dimitri and Melina Bertsekas

© 2012, 2007, 2001, 1995 Dimitri P. Bertsekas  
All rights reserved. No part of this book may be reproduced in any form  
by any electronic or mechanical means (including photocopying, recording,  
or information storage and retrieval) without permission in writing from  
the publisher.

Publisher's Cataloging-in-Publication Data

Bertsekas, Dimitri P.  
Dynamic Programming and Optimal Control  
Includes Bibliography and Index  
1. Mathematical Optimization. 2. Dynamic Programming. I. Title.  
QA402.5 .B465 2012      519.703      01-75941

ISBN-10: 1-886529-44-2, ISBN-13: 978-1-886529-44-1 (Vol. II)  
ISBN 1-886529-26-4 (Vol. I)  
ISBN 1-886529-08-6 (Two-volume set – latest editions)

# Contents

## 1. Discounted Problems – Theory

1.1. Minimization of Total Cost - Introduction . . . . .	p. 3
1.1.1. The Finite-Horizon DP Algorithm . . . . .	p. 5
1.1.2. Shorthand Notation and Monotonicity . . . . .	p. 6
1.1.3. A Preview of Infinite Horizon Results . . . . .	p. 10
1.1.4. Randomized and History-Dependent Policies . . . . .	p. 11
1.2. Discounted Problems - Bounded Cost per Stage . . . . .	p. 14
1.3. Scheduling and Multiarmed Bandit Problems . . . . .	p. 22
1.4. Discounted Continuous-Time Problems . . . . .	p. 32
1.5. The Role of Contraction Mappings . . . . .	p. 45
1.5.1. Sup-Norm Contractions . . . . .	p. 47
1.5.2. Discounted Problems - Unbounded Cost per Stage . . . . .	p. 54
1.6. General Forms of Discounted Dynamic Programming . . . . .	p. 57
1.6.1. Basic Results Under Contraction and Monotonicity . . . . .	p. 63
1.6.2. Discounted Dynamic Games . . . . .	p. 69
1.7. Notes, Sources, and Exercises . . . . .	p. 71

## 2. Discounted Problems – Computational Methods

2.1. Markovian Decision Problems . . . . .	p. 82
2.2. Value Iteration . . . . .	p. 84
2.2.1. Monotonic Error Bounds for Value Iteration . . . . .	p. 85
2.2.2. Variants of Value Iteration . . . . .	p. 92
2.2.3. Q-Learning . . . . .	p. 95
2.3. Policy Iteration . . . . .	p. 97
2.3.1. Policy Iteration for Costs . . . . .	p. 97
2.3.2. Policy Iteration for Q-Factors . . . . .	p. 102
2.3.3. Optimistic Policy Iteration . . . . .	p. 103
2.3.4. Limited Lookahead Policies and Rollout . . . . .	p. 106
2.4. Linear Programming Methods . . . . .	p. 112
2.5. Methods for General Discounted Problems . . . . .	p. 115
2.5.1. Limited Lookahead Policies and Approximations . . . . .	p. 117

2.5.2. Generalized Value Iteration . . . . .	p. 119
2.5.3. Approximate Value Iteration . . . . .	p. 120
2.5.4. Generalized Policy Iteration . . . . .	p. 123
2.5.5. Generalized Optimistic Policy Iteration . . . . .	p. 126
2.5.6. Approximate Policy Iteration . . . . .	p. 132
2.5.7. Mathematical Programming . . . . .	p. 137
2.6. Asynchronous Algorithms . . . . .	p. 138
2.6.1. Asynchronous Value Iteration . . . . .	p. 138
2.6.2. Asynchronous Policy Iteration . . . . .	p. 144
2.6.3. Policy Iteration with a Uniform Fixed Point . . . . .	p. 149
2.7. Notes, Sources, and Exercises . . . . .	p. 156
<b>3. Stochastic Shortest Path Problems</b>	
3.1. Problem Formulation . . . . .	p. 172
3.2. Main Results . . . . .	p. 175
3.3. Underlying Contraction Properties . . . . .	p. 182
3.4. Value Iteration . . . . .	p. 184
3.4.1. Conditions for Finite Termination . . . . .	p. 185
3.4.2. Asynchronous Value Iteration . . . . .	p. 188
3.5. Policy Iteration . . . . .	p. 189
3.5.1. Optimistic Policy Iteration . . . . .	p. 190
3.5.2. Approximate Policy Iteration . . . . .	p. 191
3.5.3. Policy Iteration with Improper Policies . . . . .	p. 193
3.5.4. Policy Iteration with a Uniform Fixed Point . . . . .	p. 197
3.6. Countable-State Problems . . . . .	p. 201
3.7. Notes, Sources, and Exercises . . . . .	p. 204
<b>4. Undiscounted Problems</b>	
4.1. Unbounded Costs per Stage . . . . .	p. 214
4.1.1. Main Results . . . . .	p. 216
4.1.2. Value Iteration . . . . .	p. 224
4.1.3. Other Computational Methods . . . . .	p. 230
4.2. Linear Systems and Quadratic Cost . . . . .	p. 231
4.3. Inventory Control . . . . .	p. 233
4.4. Optimal Stopping . . . . .	p. 235
4.5. Optimal Gambling Strategies . . . . .	p. 241
4.6. Continuous-Time Problems - Control of Queues . . . . .	p. 248
4.7. Nonstationary and Periodic Problems . . . . .	p. 256
4.8. Notes, Sources, and Exercises . . . . .	p. 261
<b>5. Average Cost per Stage Problems</b>	
5.1. Finite-Spaces Average Cost Models . . . . .	p. 274
5.1.1. Relation with the Discounted Cost Problem . . . . .	p. 278

- 5.1.2. Blackwell Optimal Policies . . . . . p. 284
- 5.1.3. Optimality Equations . . . . . p. 294
- 5.2. Conditions for Equal Average Cost for all Initial States . . . . . p. 298
- 5.3. Value Iteration . . . . . p. 304
  - 5.3.1. Single-Chain Value Iteration . . . . . p. 307
  - 5.3.2. Multi-Chain Value Iteration . . . . . p. 322
- 5.4. Policy Iteration . . . . . p. 329
  - 5.4.1. Single-Chain Policy Iteration . . . . . p. 329
  - 5.4.2. Multi-Chain Policy Iteration . . . . . p. 335
- 5.5. Linear Programming . . . . . p. 339
- 5.6. Infinite-Spaces Average Cost Models . . . . . p. 345
  - 5.6.1. A Sufficient Condition for Optimality . . . . . p. 353
  - 5.6.2. Finite State Space and Infinite Control Space . . . . . p. 355
  - 5.6.3. Countable States – Vanishing Discount Approach . . . . . p. 364
  - 5.6.4. Countable States – Contraction Approach . . . . . p. 367
  - 5.6.5. Linear Systems with Quadratic Cost . . . . . p. 372
- 5.7. Notes, Sources, and Exercises . . . . . p. 374

**6. Approximate Dynamic Programming - Discounted Models**

- 6.1. General Issues of Simulation-Based Cost Approximation . . . . . p. 391
  - 6.1.1. Approximation Architectures . . . . . p. 391
  - 6.1.2. Simulation-Based Approximate Policy Iteration . . . . . p. 397
  - 6.1.3. Direct and Indirect Approximation . . . . . p. 403
  - 6.1.4. Monte Carlo Simulation . . . . . p. 405
  - 6.1.5. Simplifications . . . . . p. 413
- 6.2. Direct Policy Evaluation - Gradient Methods . . . . . p. 418
- 6.3. Projected Equation Methods for Policy Evaluation . . . . . p. 423
  - 6.3.1. The Projected Bellman Equation . . . . . p. 424
  - 6.3.2. The Matrix Form of the Projected Equation . . . . . p. 428
  - 6.3.3. Simulation-Based Methods . . . . . p. 431
  - 6.3.4. LSTD, LSPE, and TD(0) Methods . . . . . p. 433
  - 6.3.5. Optimistic Versions . . . . . p. 437
  - 6.3.6. Multistep Simulation-Based Methods . . . . . p. 438
  - 6.3.7. A Synopsis . . . . . p. 447
- 6.4. Policy Iteration Issues . . . . . p. 451
  - 6.4.1. Exploration Enhancement by Geometric Sampling . . . . . p. 453
  - 6.4.2. Exploration Enhancement by Off-Policy Methods . . . . . p. 464
  - 6.4.3. Policy Oscillations - Chattering . . . . . p. 467
- 6.5. Aggregation Methods . . . . . p. 474
  - 6.5.1. Cost Approximation via the Aggregate Problem . . . . . p. 482
  - 6.5.2. Cost Approximation via the Enlarged Problem . . . . . p. 485
  - 6.5.3. Multistep Aggregation . . . . . p. 490
  - 6.5.4. Asynchronous Distributed Aggregation . . . . . p. 491
- 6.6. Q-Learning . . . . . p. 493

6.6.1. Q-Learning: A Stochastic VI Algorithm . . . . .	p. 494
6.6.2. Q-Learning and Policy Iteration . . . . .	p. 496
6.6.3. Q-Factor Approximation and Projected Equations . . . . .	p. 499
6.6.4. Q-Learning for Optimal Stopping Problems . . . . .	p. 502
6.6.5. Q-Learning and Aggregation . . . . .	p. 507
6.6.6. Finite Horizon Q-Learning . . . . .	p. 509
6.7. Notes, Sources, and Exercises . . . . .	p. 511
<b>7. Approximate Dynamic Programming - Nondiscounted Models and Generalizations</b>	
7.1. Stochastic Shortest Path Problems . . . . .	p. 532
7.2. Average Cost Problems . . . . .	p. 537
7.2.1. Approximate Policy Evaluation . . . . .	p. 537
7.2.2. Approximate Policy Iteration . . . . .	p. 546
7.2.3. Q-Learning for Average Cost Problems . . . . .	p. 548
7.3. General Problems and Monte Carlo Linear Algebra . . . . .	p. 552
7.3.1. Projected Equations . . . . .	p. 562
7.3.2. Matrix Inversion and Iterative Methods . . . . .	p. 569
7.3.3. Multistep Methods . . . . .	p. 576
7.3.4. Extension of Q-Learning for Optimal Stopping . . . . .	p. 584
7.3.5. Equation Error Methods . . . . .	p. 586
7.3.6. Oblique Projections . . . . .	p. 591
7.3.7. Generalized Aggregation . . . . .	p. 593
7.3.8. Deterministic Methods for Singular Linear Systems . . . . .	p. 597
7.3.9. Stochastic Methods for Singular Linear Systems . . . . .	p. 608
7.4. Approximation in Policy Space . . . . .	p. 620
7.4.1. The Gradient Formula . . . . .	p. 622
7.4.2. Computing the Gradient by Simulation . . . . .	p. 623
7.4.3. Essential Features for Gradient Evaluation . . . . .	p. 625
7.4.4. Approximations in Policy and Value Space . . . . .	p. 627
7.5. Notes, Sources, and Exercises . . . . .	p. 629
<b>Appendix A: Measure-Theoretic Issues in Dynamic Programming</b>	
A.1. A Two-Stage Example . . . . .	p. 641
A.2. Resolution of the Measurability Issues . . . . .	p. 646
<b>References</b> . . . . .	p. 657
<b>Index</b> . . . . .	p. 691

# CONTENTS OF VOLUME I

## 1. The Dynamic Programming Algorithm

- 1.1. Introduction
- 1.2. The Basic Problem
- 1.3. The Dynamic Programming Algorithm
- 1.4. State Augmentation and Other Reformulations
- 1.5. Some Mathematical Issues
- 1.6. Dynamic Programming and Minimax Control
- 1.7. Notes, Sources, and Exercises

## 2. Deterministic Systems and the Shortest Path Problem

- 2.1. Finite-State Systems and Shortest Paths
- 2.2. Some Shortest Path Applications
  - 2.2.1. Critical Path Analysis
  - 2.2.2. Hidden Markov Models and the Viterbi Algorithm
- 2.3. Shortest Path Algorithms
  - 2.3.1. Label Correcting Methods
  - 2.3.2. Label Correcting Variations -  $A^*$  Algorithm
  - 2.3.3. Branch-and-Bound
  - 2.3.4. Constrained and Multiobjective Problems
- 2.4. Notes, Sources, and Exercises

## 3. Deterministic Continuous-Time Optimal Control

- 3.1. Continuous-Time Optimal Control
- 3.2. The Hamilton – Jacobi – Bellman Equation
- 3.3. The Pontryagin Minimum Principle
  - 3.3.1. An Informal Derivation Using the HJB Equation
  - 3.3.2. A Derivation Based on Variational Ideas
  - 3.3.3. The Minimum Principle for Discrete-Time Problems
- 3.4. Extensions of the Minimum Principle
  - 3.4.1. Fixed Terminal State
  - 3.4.2. Free Initial State
  - 3.4.3. Free Terminal Time
  - 3.4.4. Time-Varying System and Cost
  - 3.4.5. Singular Problems
- 3.5. Notes, Sources, and Exercises

## 4. Problems with Perfect State Information

- 4.1. Linear Systems and Quadratic Cost
- 4.2. Inventory Control
- 4.3. Dynamic Portfolio Analysis
- 4.4. Optimal Stopping Problems
- 4.5. Scheduling and the Interchange Argument

- 4.6. Set-Membership Description of Uncertainty
- 4.6.1. Set-Membership Estimation
- 4.6.2. Control with Unknown-but-Bounded Disturbances
- 4.7. Notes, Sources, and Exercises
- 5. Problems with Imperfect State Information**
  - 5.1. Reduction to the Perfect Information Case
  - 5.2. Linear Systems and Quadratic Cost
  - 5.3. Minimum Variance Control of Linear Systems
  - 5.4. Sufficient Statistics and Finite-State Markov Chains
    - 5.4.1. The Conditional State Distribution
    - 5.4.2. Finite-State Systems
  - 5.5. Notes, Sources, and Exercises
- 6. Suboptimal and Adaptive Control**
  - 6.1. Certainty Equivalent and Adaptive Control
    - 6.1.1. Caution, Probing, and Dual Control
    - 6.1.2. Two-Phase Control and Identifiability
    - 6.1.3. Certainty Equivalent Control and Identifiability
    - 6.1.4. Self-Tuning Regulators
  - 6.2. Open-Loop Feedback Control
  - 6.3. Limited Lookahead Policies and Applications
    - 6.3.1. Performance Bounds for Limited Lookahead Policies
    - 6.3.2. Computational Issues in Limited Lookahead
    - 6.3.3. Problem Approximation - Enforced Decomposition
    - 6.3.4. Aggregation
    - 6.3.5. Parametric Cost-to-Go Approximation
  - 6.4. Rollout Algorithms
    - 6.4.1. Discrete Deterministic Problems
    - 6.4.2.  $Q$ -Factors Evaluated by Simulation
    - 6.4.3.  $Q$ -Factor Approximation
  - 6.5. Model Predictive Control and Related Methods
    - 6.5.1. Rolling Horizon Approximations
    - 6.5.2. Stability Issues in Model Predictive Control
    - 6.5.3. Restricted Structure Policies
  - 6.6. Additional Topics in Approximate DP
    - 6.6.1. Discretization
    - 6.6.2. Other Approximation Approaches
  - 6.7. Notes, Sources, and Exercises
- 7. Introduction to Infinite Horizon Problems**
  - 7.1. An Overview
  - 7.2. Stochastic Shortest Path Problems
  - 7.3. Discounted Problems
  - 7.4. Average Cost Problems



- 7.5. Semi-Markov Problems
- 7.6. Notes, Sources, and Exercises

### **Appendix A: Mathematical Review**

- A.1. Sets
- A.2. Euclidean Space
- A.3. Matrices
- A.4. Analysis
- A.5. Convex Sets and Functions

### **Appendix B: On Optimization Theory**

- B.1. Optimal Solutions
- B.2. Optimality Conditions
- B.3. Minimization of Quadratic Forms

### **Appendix C: On Probability Theory**

- C.1. Probability Spaces
- C.2. Random Variables
- C.3. Conditional Probability

### **Appendix D: On Finite-State Markov Chains**

- D.1. Stationary Markov Chains
- D.2. Classification of States
- D.3. Limiting Probabilities
- D.4. First Passage Times

### **Appendix E: Least-Squares Estimation and Kalman Filtering**

- E.1. Least-Squares Estimation
- E.2. Linear Least-Squares Estimation
- E.3. State Estimation – Kalman Filter
- E.4. Stability Aspects
- E.5. Gauss-Markov Estimators
- E.6. Deterministic Least-Squares Estimation

### **Appendix F: Modeling of Stochastic Linear Systems**

- F.1. Linear Systems with Stochastic Inputs
- F.2. Processes with Rational Spectrum
- F.3. The ARMAX Model

### **Appendix G: Formulating Problems of Decision Under Uncertainty**

- G.1. The Problem of Decision Under Uncertainty
- G.2. Expected Utility Theory and Risk
- G.3. Stochastic Optimal Control Problems

## ABOUT THE AUTHOR

Dimitri Bertsekas studied Mechanical and Electrical Engineering at the National Technical University of Athens, Greece, and obtained his Ph.D. in system science from the Massachusetts Institute of Technology. He has held faculty positions with the Engineering-Economic Systems Dept., Stanford University, and the Electrical Engineering Dept. of the University of Illinois, Urbana. Since 1979 he has been teaching at the Electrical Engineering and Computer Science Department of the Massachusetts Institute of Technology (M.I.T.), where he is currently McAfee Professor of Engineering.

His research spans several fields, including optimization, control, large-scale computation, and data communication networks, and is closely tied to his teaching and book authoring activities. He has written numerous research papers, and thirteen books, several of which are used as textbooks in MIT classes. He consults regularly with private industry and has held editorial positions in several journals.

Professor Bertsekas was awarded the INFORMS 1997 Prize for Research Excellence in the Interface Between Operations Research and Computer Science for his book “Neuro-Dynamic Programming” (co-authored with John Tsitsiklis), the 2000 Greek National Award for Operations Research, the 2001 ACC John R. Ragazzini Education Award, and the 2009 INFORMS Expository Writing Award. In 2001, he was elected to the United States National Academy of Engineering.

# *Preface*

This two-volume book is based on a first-year graduate course on dynamic programming and optimal control that I have taught for over twenty years at Stanford University, the University of Illinois, and the Massachusetts Institute of Technology. The course has been typically attended by students from engineering, operations research, economics, and applied mathematics. Accordingly, a principal objective of the book has been to provide a unified treatment of the subject, suitable for a broad audience. In particular, problems with a continuous character, such as stochastic control problems, popular in modern control theory, are simultaneously treated with problems with a discrete character, such as Markovian decision problems, popular in operations research. Furthermore, many applications and examples, drawn from a broad variety of fields, are discussed.

The book may be viewed as a greatly expanded and pedagogically improved version of my 1987 book “Dynamic Programming: Deterministic and Stochastic Models,” published by Prentice-Hall. I have included much new material on deterministic and stochastic shortest path problems, as well as a new chapter on continuous-time optimal control problems and the Pontryagin Maximum Principle, developed from a dynamic programming viewpoint. I have also added a fairly extensive exposition of simulation-based approximation techniques for dynamic programming. These techniques, which are often referred to as “neuro-dynamic programming” or “reinforcement learning,” represent a breakthrough in the practical application of dynamic programming to complex problems that involve the dual curse of large dimension and lack of an accurate mathematical model. Other material was also augmented, substantially modified, and updated.

With the new material, however, the book grew so much in size that it became necessary to divide it into two volumes: one on finite horizon, and the other on infinite horizon problems. This division was not only natural in terms of size, but also in terms of style and orientation. The first volume is more oriented towards modeling, and the second is more oriented towards mathematical analysis and computation. I have included in the first volume a final chapter that provides an introductory treatment of infinite horizon problems. The purpose is to make the first volume self-

contained for instructors who wish to cover a modest amount of infinite horizon material in a course that is primarily oriented towards modeling, conceptualization, and finite horizon problems,

Many topics in the book are relatively independent of the others. For example Chapter 2 of Vol. I on shortest path problems can be skipped without loss of continuity, and the same is true for Chapter 3 of Vol. I, which deals with continuous-time optimal control. As a result, the book can be used to teach several different types of courses.

- (a) A two-semester course that covers both volumes.
- (b) A one-semester course primarily focused on finite horizon problems that covers most of the first volume.
- (c) A one-semester course focused on stochastic optimal control that covers Chapters 1, 4, 5, and 6 of Vol. I, and Chapters 1, 2, and 4 of Vol. II.
- (d) A one-semester course that covers Chapter 1, about 50% of Chapters 2 through 6 of Vol. I, and about 70% of Chapters 1, 2, and 4 of Vol. II. This is the course I usually teach at MIT.
- (e) A one-quarter engineering course that covers the first three chapters and parts of Chapters 4 through 6 of Vol. I.
- (f) A one-quarter mathematically oriented course focused on infinite horizon problems that covers Vol. II.

The mathematical prerequisite for the text is knowledge of advanced calculus, introductory probability theory, and matrix-vector algebra. A summary of this material is provided in the appendixes. Naturally, prior exposure to dynamic system theory, control, optimization, or operations research will be helpful to the reader, but based on my experience, the material given here is reasonably self-contained.

The book contains a large number of exercises, and the serious reader will benefit greatly by going through them. Solutions to all exercises are compiled in a manual that is available to instructors from the author. Many thanks are due to the several people who spent long hours contributing to this manual, particularly Steven Shreve, Eric Loiederman, Lakis Polymenakos, and Cynara Wu.

Dynamic programming is a conceptually simple technique that can be adequately explained using elementary analysis. Yet a mathematically rigorous treatment of general dynamic programming requires the complicated machinery of measure-theoretic probability. My choice has been to bypass the complicated mathematics by developing the subject in generality, while claiming rigor only when the underlying probability spaces are countable. A mathematically rigorous treatment of the subject is carried out in my monograph “Stochastic Optimal Control: The Discrete Time

Case,” Academic Press, 1978,† coauthored by Steven Shreve. This monograph complements the present text and provides a solid foundation for the subjects developed somewhat informally here.

Finally, I am thankful to a number of individuals and institutions for their contributions to the book. My understanding of the subject was sharpened while I worked with Steven Shreve on our 1978 monograph. My interaction and collaboration with John Tsitsiklis on stochastic shortest paths and approximate dynamic programming have been most valuable. Michael Caramanis, Emmanuel Fernandez-Gaucherand, Pierre Humblet, Lennart Ljung, and John Tsitsiklis taught from versions of the book, and contributed several substantive comments and homework problems. A number of colleagues offered valuable insights and information, particularly David Castanon, Eugene Feinberg, and Krishna Pattipati. NSF provided research support. Prentice-Hall graciously allowed the use of material from my 1987 book. Teaching and interacting with the students at MIT have kept up my interest and excitement for the subject.

Dimitri P. Bertsekas  
bertsekas@lids.mit.edu  
<http://web.mit.edu/dimitrib/www/home.html>

---

† Note added in the 2nd edition: This monograph was republished by Athena Scientific in 1996.

## *Preface to the Second Edition*

This second edition of Vol. II should be viewed as a relatively minor revision of the original. The coverage was expanded in a few areas as follows:

- (a) In Chapter 1, material was added on variants of the policy iteration method.
- (b) In Chapter 2, the material on neuro-dynamic programming methods was updated and expanded to reflect some recent developments.
- (c) In Chapter 4, material was added on some new value iteration methods.
- (d) In Chapter 5, the material on semi-Markov problems was revised, with a significant portion simplified and shifted to Volume I.

There are also a few miscellaneous additions and improvements scattered throughout the text. Finally, a new internet-based feature was added to the book, which extends its scope and coverage. Many of the theoretical exercises have been solved in detail and their solutions have been posted in the book's www page

<http://www.athenasc.com/dpbook.html>

These exercises have been marked with the symbol

I would like to express my thanks to the many colleagues who contributed suggestions for improvement of the second edition.

Dimitri P. Bertsekas  
bertsekas@lids.mit.edu  
<http://web.mit.edu/dimitrib/www/home.html>

June, 2001

## *Preface to the Third Edition*

This is a major revision of the 2nd edition, and contains a substantial amount of new material, as well as a major reorganization of old material. The length of the text has increased by more than 50%, and more than half of the old material has been restructured and/or revised. Most of the added material is in four areas.

- (a) The coverage of the average cost problem of Chapter 4 has greatly increased in scope and depth. In particular, there is now a full analysis of multi-chain problems, as well as a more extensive analysis of infinite-spaces problems (Section 4.6).
- (b) The material on approximate dynamic programming has been collected in Chapter 6. It has been greatly expanded to include new research, thereby supplementing the 1996 book “Neuro-Dynamic Programming.”
- (c) Contraction mappings and their role in various analyses have been highlighted in new material on infinite state space problems (Sections 1.4, 2.5, and 4.6), and in their use in the approximate dynamic programming material of Chapter 6.
- (d) The mathematical measure-theoretic issues that must be addressed for a rigorous theory of stochastic dynamic programming have been illustrated and summarized in an appendix for the benefit of the mathematically oriented reader.

Also some exercises were added and a few sections were revised while preserving their essential content.

I would like to express my thanks to many colleagues who contributed valuable comments. I am particularly thankful to Ciamac Moallemi, Steven Shreve, John Tsitsiklis, and Ben Van Roy, who reviewed some of the new material and each contributed several substantial suggestions. I wish to thank especially Janey Yu who read with great care and keen eye large parts of the book, contributed important analysis and many incisive, substantive comments, and also collaborated with me in research that was included in Chapter 6.

Dimitri P. Bertsekas  
<http://web.mit.edu/dimitrib/www/home.html>  
Fall 2006

## *Preface to the Fourth Edition*

This is a major revision of Vol. II, and contains a substantial amount of new material, as well as a reorganization of old material. The length has increased by more than 60% from the third edition, and most of the old material has been restructured and/or revised. Volume II now numbers more than 700 pages and is larger in size than Vol. I. It can arguably be viewed as a new book!

Approximate DP has become the central focal point of Vol. II, and occupies more than half of the book (the last two chapters, and large parts of Chapters 1-3). Thus one may also view Vol. II as a followup to my 1996 book “Neuro-Dynamic Programming” (coauthored with John Tsitsiklis). The present book focuses to a great extent on new research that became available after 1996. On the other hand, the textbook style of the book has been preserved, and some mathematically demanding material has been explained at an intuitive or informal level, while referring to the journal literature or the Neuro-Dynamic Programming book for a more mathematical treatment.

In the process of expansion and reorganization, the design of the book became more modular and suitable for classroom use. The core material, which can be covered in about a third to a half of one semester is Chapter 1 (except for the application-specific Sections 1.3 and 1.4), Chapter 2, and Chapter 6, which are self-contained when taken together. This material focuses on discounted problems, and may be supplemented by parts of Chapter 3 and Section 7.1 on stochastic shortest path problems. Indeed, this comprises half of what the author covers in his MIT class (the remaining half comes from Volume I, including Chapter 6 of that volume that deals with finite horizon approximate DP). The material on average cost problems, given in Chapter 5, and Sections 7.2 and 7.4, and the advanced material on positive and negative DP models (Chapter 4), and Monte Carlo linear algebra (Section 7.3) are terminal subjects that may be covered at the instructor’s discretion.

As our focus shifted, we have placed increased emphasis on new or recent research in approximate DP and simulation-based methods, as well as on asynchronous iterative methods, in view of the central role of simulation, which is by nature asynchronous. A lot of this material is an outgrowth of the author’s research in the six years since the third edition was published. Some of the highlights, in the order appearing in the book, are:

- (a) Computational methods for generalized discounted DP (Sections 2.5 and 2.6), including the asynchronous optimistic policy iteration ma-



- terial of Section 2.6.3, and its application to game and minimax problems, constrained policy iteration, and Q-learning.
- (b) Policy iteration methods (including asynchronous optimistic versions) for stochastic shortest path problems that involve improper policies (Section 3.4).
  - (c) Extensive new material on various simulation-based, approximate value and policy iteration methods in Sections 6.3-6.6 (projected equation methods, aggregation methods, and Q-learning).
  - (d) New simulation techniques for multistep methods, such as geometric and free-form sampling (Sections 6.4.1 and 7.3.3).
  - (e) Extensive new material on Monte Carlo linear algebra in Section 7.3 (primarily the simulation-based and approximate solution of large systems of linear equations), which extends the DP methodology of approximate policy evaluation.

Much of the research in (a)-(d) is based on my work with Janey (Huizhen) Yu, while the research in (e) is based on my work with Janey Yu and Mengdi Wang. My collaboration with Janey and Mengdi has been very productive and is greatly appreciated. The reader is referred to our joint and individual papers, which describe more fully our research, including material that could not be covered in this book.

I want to express my appreciation to colleagues and collaborators in approximate DP research, who contributed to the book in various ways, particularly Vivek Borkar, Angelia Nedić, and Ben Van Roy. A special thanks goes to John Tsitsiklis, with whom I have interacted extensively through collaboration and sharing of ideas on DP and asynchronous algorithms for more than 30 years. I also wish to acknowledge helpful interactions with many colleagues, including Vivek Farias, Eugene Feinberg, Warren Powell, Martin Puterman, Uriel Rothblum, and Bruno Scherrer. Finally, I want to thank the many students in my DP classes of the last decade, who patiently labored with a textbook under development, and contributed their ideas and experiences through their research projects from a broad variety of application fields.

Dimitri P. Bertsekas  
<http://web.mit.edu/dimitrib/www/home.html>  
Spring 2012