

*Dynamic Programming
and Optimal Control*
Volume I

Excerpts For ASU Course Use

Please do not Distribute for Other Uses

Dimitri P. Bertsekas

Massachusetts Institute of Technology

WWW site for book information and orders

<http://www.athenasc.com>



Athena Scientific, Belmont, Massachusetts

About this Material

I provide selected sections from Chapters 1 and 5 of Vol. I of the two-volume book *Dynamic Programming and Optimal Control* (4th edition, 2017) for use in Arizona State University classes. The book was developed as a textbook for a first-year graduate course on dynamic programming and optimal control that I have taught since the mid 70s at Stanford University, the University of Illinois, and the Massachusetts Institute of Technology.

Course material from several offerings of my class can be found at my website and at the MIT Open CourseWare (OCW) site

<https://ocw.mit.edu/index.htm>

including slides, selected exercise solutions, and other material.

Links to video lectures on approximate DP and related topics may be found at my website, which also contains my research papers on the subject. The full book may be ordered from the publishing company or Amazon.com.

Dimitri P. Bertsekas

<http://web.mit.edu/dimitrib/www/home.html>

Fall 2018

Contents

1. The Dynamic Programming Algorithm

1.1. Introduction	p. 2
1.1.1. General Structure of Finite Horizon Optimal Control Problems	p. 4
1.1.2. Discrete-State and Finite-State Problems	p. 7
1.2. The Basic Problem	p. 14
1.3. The Dynamic Programming Algorithm	p. 20
1.4. State Augmentation and Other Reformulations	p. 37
1.5. Some Mathematical Issues	p. 44
1.6. Dynamic Programming and Minimax Control	p. 49
1.7. Notes, Sources, and Exercises	p. 53

2. Deterministic Systems and the Shortest Path Problem

2.1. Finite-State Systems and Shortest Paths	p. 69
2.2. Some Shortest Path Applications	p. 72
2.2.1. Critical Path Analysis	p. 72
2.2.2. Hidden Markov Models and the Viterbi Algorithm	p. 74
2.3. Shortest Path Algorithms	p. 81
2.3.1. Label Correcting Methods	p. 82
2.3.2. Label Correcting Variations - A^* Algorithm	p. 91
2.3.3. Branch-and-Bound	p. 92
2.3.4. Constrained and Multiobjective Problems	p. 95
2.4. Notes, Sources, and Exercises	p. 101

3. Problems with Perfect State Information

3.1. Linear Systems and Quadratic Cost	p. 110
3.2. Inventory Control	p. 125
3.3. Dynamic Portfolio Analysis	p. 134
3.4. Optimal Stopping Problems	p. 140

3.5. Scheduling and the Interchange Argument	p. 150
3.6. Set-Membership Description of Uncertainty	p. 154
3.6.1. Set-Membership Estimation	p. 155
3.6.2. Control within a Target Tube	p. 161
3.7. Notes, Sources, and Exercises	p. 165
4. Problems with Imperfect State Information	
4.1. Reduction to the Perfect Information Case	p. 184
4.2. Linear Systems and Quadratic Cost	p. 195
4.3. Sufficient Statistics	p. 202
4.3.1. The Conditional State Distribution	p. 204
4.3.2. Finite-State Systems	p. 209
4.4. Notes, Sources, and Exercises	p. 221
5. Introduction to Infinite Horizon Problems	
5.1. An Overview	p. 232
5.2. Stochastic Shortest Path Problems	p. 236
5.3. Computational Methods	p. 245
5.3.1. Value Iteration	p. 245
5.3.2. Policy Iteration	p. 246
5.3.3. Linear Programming	p. 248
5.4. Discounted Problems	p. 249
5.5. Average Cost per Stage Problems	p. 253
5.6. Semi-Markov Problems	p. 267
5.7. Notes, Sources, and Exercises	p. 277
6. Approximate Dynamic Programming	
6.1. Cost Approximation and Limited Lookahead	p. 296
6.1.1. Error Bounds and Cost Improvement	p. 300
6.1.2. Computation of Suboptimal Policies - Stochastic Programming	p. 304
6.2. Problem Approximation	p. 307
6.2.1. Enforced Decomposition	p. 307
6.2.2. Probabilistic Approximation - Certainty Equivalent Control	p. 316
6.2.3. Aggregation	p. 321
6.3. Parametric Cost Approximation	p. 327
6.3.1. Feature-Based Architectures and Neural Networks	p. 327
6.3.2. Sequential Dynamic Programming Approximation	p. 337
6.3.3. Q -Factor Parametric Approximation	p. 339
6.3.4. Parametric Approximation in Infinite Horizon Problems	p. 342

6.3.5. Computer Chess	p. 345
6.4. On-Line Approximation and Optimization	p. 352
6.4.1. Rollout Algorithms	p. 352
6.4.2. Rollout for Discrete Deterministic Problems	p. 360
6.4.3. Model Predictive Control	p. 376
6.4.4. Open-Loop Feedback Control	p. 385
6.5. Simulation-Based Cost-to-go Approximation	p. 389
6.5.1. Stochastic Rollout and Monte Carlo Tree Search	p. 390
6.5.2. Variance Reduction in Rollout	p. 392
6.6. Approximation in Policy Space	p. 395
6.7. Adaptive Control	p. 397
6.8. Discretization Issues	p. 405
6.9. Notes, Sources, and Exercises	p. 408

7. Deterministic Continuous-Time Optimal Control

7.1. Continuous-Time Optimal Control	p. 426
7.2. The Hamilton-Jacobi-Bellman Equation	p. 429
7.3. The Pontryagin Minimum Principle	p. 435
7.3.1. An Informal Derivation Using the HJB Equation	p. 435
7.3.2. A Derivation Based on Variational Ideas	p. 445
7.3.3. Minimum Principle for Discrete-Time Problems	p. 449
7.4. Extensions of the Minimum Principle	p. 451
7.4.1. Fixed Terminal State	p. 451
7.4.2. Free Initial State	p. 455
7.4.3. Free Terminal Time	p. 455
7.4.4. Time-Varying System and Cost	p. 459
7.4.5. Singular Problems	p. 459
7.5. Notes, Sources, and Exercises	p. 461

Appendix A: Mathematical Review

A.1. Sets	p. 467
A.2. Euclidean Space	p. 468
A.3. Matrices	p. 469
A.4. Analysis	p. 473
A.5. Convex Sets and Functions	p. 475

Appendix B: On Optimization Theory

B.1. Optimal Solutions	p. 476
B.2. Optimality Conditions	p. 477
B.3. Minimization of Quadratic Forms	p. 479

Appendix C: On Probability Theory

C.1. Probability Spaces	p. 480
C.2. Random Variables	p. 481
C.3. Conditional Probability	p. 482

Appendix D: On Finite-State Markov Chains

D.1. Stationary Markov Chains	p. 485
D.2. Classification of States	p. 486
D.3. Limiting Probabilities	p. 487
D.4. First Passage Times	p. 488

Appendix E: Least Squares Estimation and Kalman Filtering

E.1. Least-Squares Estimation	p. 489
E.2. Linear Least-Squares Estimation	p. 491
E.3. State Estimation – Kalman Filter	p. 499
E.4. Stability Aspects	p. 504
E.5. Gauss-Markov Estimators	p. 507
E.6. Deterministic Least-Squares Estimation	p. 509

Appendix F: Formulating Problems of Decision Under Uncertainty

F.1. The Problem of Decision Under Uncertainty	p. 511
F.2. Expected Utility Theory and Risk	p. 515
F.3. Stochastic Optimal Control Problems	p. 528

References	p. 533
-----------------------------	---------------

Index	p. 551
------------------------	---------------

CONTENTS OF VOLUME II

1. Discounted Problems – Theory

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

2. Discounted Problems – Computational Methods

- 2.1. Markovian Decision Problems
- 2.2. Value Iteration
 - 2.2.1. Monotonic Error Bounds for Value Iteration
 - 2.2.2. Variants of Value Iteration
 - 2.2.3. Q-Learning
- 2.3. Policy Iteration
 - 2.3.1. Policy Iteration for Costs
 - 2.3.2. Policy Iteration for Q-Factors
 - 2.3.3. Optimistic Policy Iteration
 - 2.3.4. Limited Lookahead Policies and Rollout
- 2.4. Linear Programming Methods
- 2.5. Methods for General Discounted Problems
 - 2.5.1. Limited Lookahead Policies and Approximations
 - 2.5.2. Generalized Value Iteration
 - 2.5.3. Approximate Value Iteration
 - 2.5.4. Generalized Policy Iteration
 - 2.5.5. Generalized Optimistic Policy Iteration
 - 2.5.6. Approximate Policy Iteration
 - 2.5.7. Mathematical Programming
- 2.6. Asynchronous Algorithms
 - 2.6.1. Asynchronous Value Iteration

- 2.6.2. Asynchronous Policy Iteration
- 2.6.3. Policy Iteration with a Uniform Fixed Point
- 2.7. Notes, Sources, and Exercises

3. Stochastic Shortest Path Problems

- 3.1. Problem Formulation
- 3.2. Main Results
- 3.3. Underlying Contraction Properties
- 3.4. Value Iteration
 - 3.4.1. Conditions for Finite Termination
 - 3.4.2. Asynchronous Value Iteration
- 3.5. Policy Iteration
 - 3.5.1. Optimistic Policy Iteration
 - 3.5.2. Approximate Policy Iteration
 - 3.5.3. Policy Iteration with Improper Policies
 - 3.5.4. Policy Iteration with a Uniform Fixed Point
- 3.6. Countable-State Problems
- 3.7. Notes, Sources, and Exercises

4. Undiscounted Problems

- 4.1. Unbounded Costs per Stage
 - 4.1.1. Main Results
 - 4.1.2. Value Iteration
 - 4.1.3. Other Computational Methods
- 4.2. Linear Systems and Quadratic Cost
- 4.3. Inventory Control
- 4.4. Optimal Stopping
- 4.5. Optimal Gambling Strategies
- 4.6. Continuous-Time Problems - Control of Queues
- 4.7. Nonstationary and Periodic Problems
- 4.8. Notes, Sources, and Exercises

5. Average Cost per Stage Problems

- 5.1. Finite-Spaces Average Cost Models
 - 5.1.1. Relation with the Discounted Cost Problem
 - 5.1.2. Blackwell Optimal Policies
 - 5.1.3. Optimality Equations
- 5.2. Conditions for Equal Average Cost for all Initial States
- 5.3. Value Iteration
 - 5.3.1. Single-Chain Value Iteration
 - 5.3.2. Multi-Chain Value Iteration
- 5.4. Policy Iteration
 - 5.4.1. Single-Chain Policy Iteration

- 5.4.2. Multi-Chain Policy Iteration
- 5.5. Linear Programming
- 5.6. Infinite-Spaces Average Cost Models
 - 5.6.1. A Sufficient Condition for Optimality
 - 5.6.2. Finite State Space and Infinite Control Space
 - 5.6.3. Countable States – Vanishing Discount Approach
 - 5.6.4. Countable States – Contraction Approach
 - 5.6.5. Linear Systems with Quadratic Cost
- 5.7. Notes, Sources, and Exercises

6. Approximate Dynamic Programming - Discounted Models

- 6.1. General Issues of Simulation-Based Cost Approximation
 - 6.1.1. Approximation Architectures
 - 6.1.2. Simulation-Based Approximate Policy Iteration
 - 6.1.3. Direct and Indirect Approximation
 - 6.1.4. Monte Carlo Simulation
 - 6.1.5. Simplifications
- 6.2. Direct Policy Evaluation - Gradient Methods
- 6.3. Projected Equation Methods for Policy Evaluation
 - 6.3.1. The Projected Bellman Equation
 - 6.3.2. The Matrix Form of the Projected Equation
 - 6.3.3. Simulation-Based Methods
 - 6.3.4. LSTD, LSPE, and TD(0) Methods
 - 6.3.5. Optimistic Versions
 - 6.3.6. Multistep Simulation-Based Methods
 - 6.3.7. A Synopsis
- 6.4. Policy Iteration Issues
 - 6.4.1. Exploration Enhancement by Geometric Sampling
 - 6.4.2. Exploration Enhancement by Off-Policy Methods
 - 6.4.3. Policy Oscillations - Chattering
- 6.5. Aggregation Methods
 - 6.5.1. Cost Approximation via the Aggregate Problem
 - 6.5.2. Cost Approximation via the Enlarged Problem
 - 6.5.3. Multistep Aggregation
 - 6.5.4. Asynchronous Distributed Aggregation
- 6.6. Q-Learning
 - 6.6.1. Q-Learning: A Stochastic VI Algorithm
 - 6.6.2. Q-Learning and Policy Iteration
 - 6.6.3. Q-Factor Approximation and Projected Equations
 - 6.6.4. Q-Learning for Optimal Stopping Problems
 - 6.6.5. Q-Learning and Aggregation
 - 6.6.6. Finite Horizon Q-Learning
- 6.7. Notes, Sources, and Exercises

7. Approximate Dynamic Programming - Nondiscounted Models and Generalizations

- 7.1. Stochastic Shortest Path Problems
- 7.2. Average Cost Problems
 - 7.2.1. Approximate Policy Evaluation
 - 7.2.2. Approximate Policy Iteration
 - 7.2.3. Q-Learning for Average Cost Problems
- 7.3. General Problems and Monte Carlo Linear Algebra
 - 7.3.1. Projected Equations
 - 7.3.2. Matrix Inversion and Iterative Methods
 - 7.3.3. Multistep Methods
 - 7.3.4. Extension of Q-Learning for Optimal Stopping
 - 7.3.5. Equation Error Methods
 - 7.3.6. Oblique Projections
 - 7.3.7. Generalized Aggregation
 - 7.3.8. Deterministic Methods for Singular Linear Systems
 - 7.3.9. Stochastic Methods for Singular Linear Systems
- 7.4. Approximation in Policy Space
 - 7.4.1. The Gradient Formula
 - 7.4.2. Computing the Gradient by Simulation
 - 7.4.3. Essential Features for Gradient Evaluation
 - 7.4.4. Approximations in Policy and Value Space
- 7.5. Notes, Sources, and Exercises

Appendix A: Measure-Theoretic Issues in Dynamic Programming

- A.1. A Two-Stage Example
- A.2. Resolution of the Measurability Issues

References**Index**

ATHENA SCIENTIFIC
RELATED TITLES OF INTEREST

1. Dynamic Programming and Optimal Control, Two-Volume Set, by Dimitri P. Bertsekas, 2017, ISBN 1-886529-08-6, 1270 pages
2. Neuro-Dynamic Programming, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1996, ISBN 1-886529-10-8, 512 pages
3. Abstract Dynamic Programming, 2nd Ed., by Dimitri P. Bertsekas, 2018, ISBN 78-1-886529-46-5, 360 pages
4. Stochastic Optimal Control: The Discrete-Time Case, by Dimitri P. Bertsekas and Steven E. Shreve, 1996, ISBN 1-886529-03-5, 330 pages
5. Introduction to Probability, 2nd Edition, by Dimitri P. Bertsekas and John N. Tsitsiklis, 2008, ISBN 978-1-886529-23-6, 544 pages
6. Nonlinear Programming, 3rd Edition, by Dimitri P. Bertsekas, 2016, ISBN 1-886529-05-1, 880 pages

The Dynamic Programming Algorithm

Contents

1.1. Introduction	p. 2
1.1.1. General Structure of Finite Horizon Optimal Control Problems	p. 4
1.1.2. Discrete-State and Finite-State Problems	p. 7
1.2. The Basic Problem	p. 14
1.3. The Dynamic Programming Algorithm	p. 20
1.4. State Augmentation and Other Reformulations	p. 37
1.5. Some Mathematical Issues	p. 44
1.6. Dynamic Programming and Minimax Control	p. 49
1.7. Notes, Sources, and Exercises	p. 53

**Life can only be understood going backwards,
but it must be lived going forwards.
Kierkegaard**

1.1 INTRODUCTION

This book deals with situations where decisions are made in stages. The outcome of each decision may not be fully predictable but can be anticipated to some extent before the next decision is made. The objective is to minimize a certain cost over a given number of stages – a mathematical expression of what is considered an undesirable outcome.

A key aspect of such situations is that decisions cannot be viewed in isolation since one must balance the desire for low present cost with the undesirability of high future costs. The dynamic programming (DP) technique captures this tradeoff. At each stage, it ranks decisions based on the sum of the present cost and the expected future cost, assuming optimal decision making for subsequent stages.

There is a very broad variety of practical problems that can be treated by DP. In this book, we try to keep the main ideas uncluttered by irrelevant assumptions on problem structure. To this end, we formulate a broadly applicable model of stochastic optimal control of a dynamic system with perfect state observations over a finite number of stages (a finite horizon). This model will be the starting point for our development, and will occupy us for the first four of the seven chapters of this volume (which may be viewed as Part I of the book). The last three chapters (which may be viewed as Part II of the book) deal with related topics, and are terminal chapters for the purposes of this volume. In fact each of these chapters may be attempted by some readers immediately after Chapter 1, with relatively little loss of continuity. In summary, the seven chapters are structured as follows (see Fig. 1.1.1):

- (1) The first chapter deals with the formulation of a general optimal control problem over a finite horizon, it demonstrates its broad applicability in deterministic and stochastic settings, and develops the DP algorithm for its solution.
- (2) The second chapter considers the deterministic finite-state case of the problem. It explores the connections with the classical shortest path problem, and the special algorithms that this connection brings to bear.

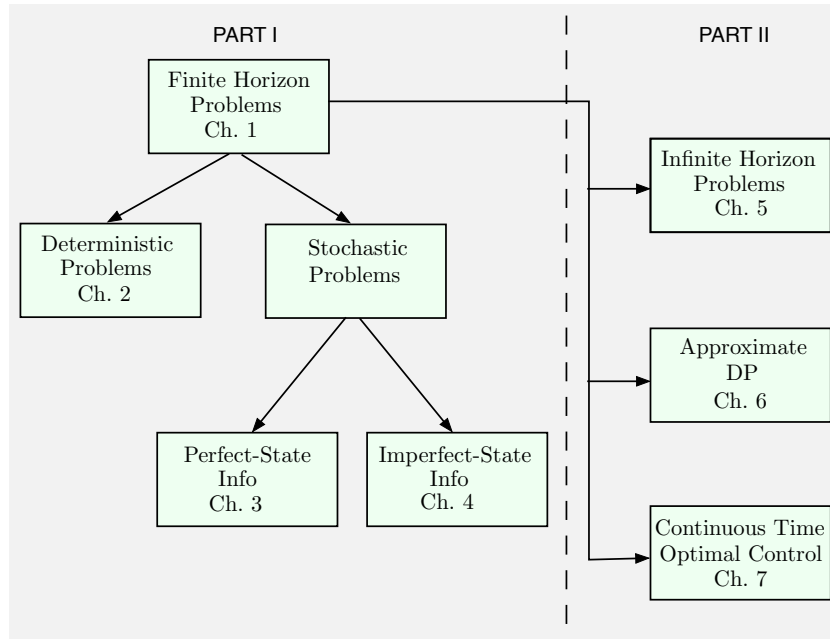


Figure 1.1.1 Illustration of the structure of the seven chapters of the present volume.

- (3) The third chapter considers the stochastic general-state case of the problem, and illustrates various aspects of the solution process by means of some important applications.
- (4) The fourth chapter also considers the stochastic general-state case, but contrary to the third chapter, it considers the situation where the exact state of the system is not observed perfectly by the decision maker/controller. This is a much harder problem, but conceptually it is closely related to the case of perfect state observation. Much of the chapter is devoted to explaining this important conceptual connection.
- (5) The fifth chapter is an introduction to the theory of infinite horizon problems. It focuses on the relatively easy but important case of finite state problems. Volume II considers infinite horizon problems in greater generality and depth.
- (6) The sixth chapter considers approximations to the exact DP solution method. This is a subject of great importance in practice, with a rich algorithmic methodology and very broad applications. We focus here primarily on finite horizon problems, so that this chapter can be read independently of Chapter 5. However, much of the discussion extends to infinite horizon problems, and on occasion we pause to indicate the

nature of the extension. We will consider approximate DP for infinite horizon problems in greater detail in Vol. II.

- (7) The seventh chapter is a terminal chapter on deterministic optimal control in continuous space and time. It may be skipped without loss of continuity. Alternatively, it may be read immediately following Chapter 1. Among others, we emphasize here the methodological connections with DP and the analog of the DP algorithm in continuous time, which is the Hamilton-Jacobi-Bellman equation.

1.1.1 General Structure of Finite Horizon Optimal Control Problems

Our finite horizon model has two principal features: (1) a *discrete-time dynamic system*, and (2) a *cost function that is additive over time*. The dynamic system expresses the evolution of some variables, the system's "state," under the influence of decisions made at discrete instances of time. The system has the form

$$x_{k+1} = f_k(x_k, u_k, w_k), \quad k = 0, 1, \dots, N-1,$$

where

k indexes discrete time,

x_k is the state of the system and summarizes past information that is relevant for future optimization,

u_k is the control or decision variable to be selected at time k ,

w_k is a random parameter (also called disturbance or noise depending on the context),

N is the horizon or number of times control is applied,

and f_k is a function that describes the system and in particular the mechanism by which the state is updated.

The cost function is additive in the sense that the cost incurred at time k , denoted by $g_k(x_k, u_k, w_k)$, accumulates over time. The total cost is

$$g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k),$$

where $g_N(x_N)$ is a terminal cost incurred at the end of the process. However, because of the presence of w_k , the cost is generally a random variable and cannot be meaningfully optimized. We therefore formulate the problem as an optimization of the *expected cost*

$$E \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right\},$$

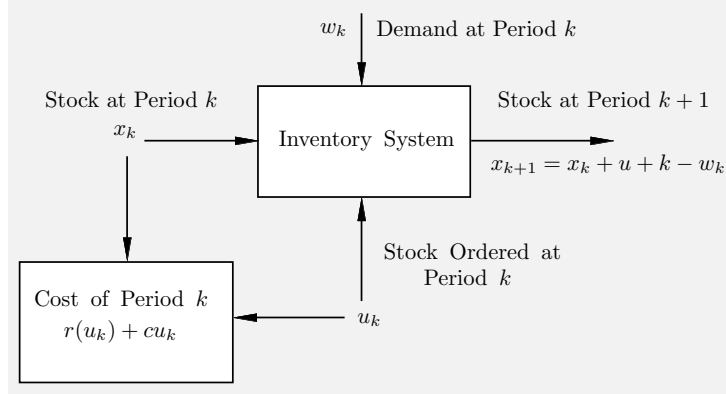


Figure 1.1.2 Inventory control example. At period k , the current stock (state) x_k , the stock ordered (control) u_k , and the demand (random disturbance) w_k determine the cost $r(x_k) + cu_k$ of period k and the stock $x_{k+1} = x_k + u_k - w_k$ at the next period.

where the expectation is with respect to the joint distribution of the random variables involved. The optimization is over the controls u_0, u_1, \dots, u_{N-1} , but some qualification is needed here: each control u_k is selected with some knowledge of the current state x_k (either its exact value or some other related information).

A more precise definition of the terminology just used will be given shortly. We first provide some orientation by means of examples.

Example 1.1.1 (Inventory Control - Open-Loop and Closed-Loop Optimization)

Consider a problem of ordering a quantity of a certain item at each of N periods so as to (roughly) meet a stochastic demand, while minimizing the incurred expected cost. Let us denote

x_k stock available at the beginning of the k th period,

u_k stock ordered (and immediately delivered) at the beginning of the k th period,

w_k demand during the k th period with given probability distribution.

We assume that w_0, w_1, \dots, w_{N-1} are independent random variables, and that excess demand is backlogged and filled as soon as additional inventory becomes available. Thus, stock evolves according to the discrete-time equation

$$x_{k+1} = x_k + u_k - w_k,$$

where negative stock corresponds to backlogged demand (see Fig. 1.1.2).

The cost incurred in period k consists of two components:

- (a) A cost $r(x_k)$ representing a penalty for either positive stock x_k (holding cost for excess inventory) or negative stock x_k (shortage cost for unfilled demand).
- (b) The purchasing cost cu_k , where c is cost per unit ordered.

There is also a terminal cost $R(x_N)$ for being left with inventory x_N at the end of N periods. Thus, the total cost over N periods is

$$E \left\{ R(x_N) + \sum_{k=0}^{N-1} (r(x_k) + cu_k) \right\}.$$

We want to minimize this cost by proper choice of the orders u_0, \dots, u_{N-1} , subject to the natural constraint $u_k \geq 0$ for all k .

At this point we need to distinguish between *closed-loop* and *open-loop* minimization of the cost. In open-loop minimization we select all orders u_0, \dots, u_{N-1} at once at time 0, without waiting to see the subsequent demand levels. In closed-loop minimization we postpone placing the order u_k until the last possible moment (time k) when the current stock x_k will be known. The idea is that since there is no penalty for delaying the order u_k up to time k , we can take advantage of information that becomes available between times 0 and k (the demand and stock level in past periods).

Closed-loop optimization is of central importance in DP and is the type of optimization that we will consider almost exclusively in this book. Thus, in our basic formulation, decisions are made in stages while gathering information between stages that will be used to enhance the quality of the decisions. The effect of this on the structure of the resulting optimization problem is quite profound. In particular, in closed-loop inventory optimization we are not interested in finding optimal numerical values of the orders but rather we want to find an *optimal rule for selecting at each period k an order u_k for each possible value of stock x_k that can conceivably occur*. This is an “action versus strategy” distinction.

Mathematically, in closed-loop inventory optimization, we want to find a sequence of functions μ_k , $k = 0, \dots, N-1$, mapping stock x_k into order u_k so as to minimize the expected cost. The meaning of μ_k is that, for each k and each possible value of x_k ,

$$\mu_k(x_k) = \text{amount that should be ordered at time } k \text{ if the stock is } x_k.$$

The sequence $\pi = \{\mu_0, \dots, \mu_{N-1}\}$ will be referred to as a *policy* or *control law*. For each π , the corresponding cost for a fixed initial stock x_0 is

$$J_\pi(x_0) = E \left\{ R(x_N) + \sum_{k=0}^{N-1} (r(x_k) + c\mu_k(x_k)) \right\},$$

and we want to minimize $J_\pi(x_0)$ for a given x_0 over all π that satisfy the constraints of the problem. This is a typical DP problem. We will analyze this problem in various forms in subsequent sections. For example, we will

show in Section 3.2 that for a reasonable choice of the cost function, the optimal ordering policy is of the form

$$\mu_k(x_k) = \begin{cases} S_k - x_k & \text{if } x_k < S_k, \\ 0 & \text{otherwise,} \end{cases}$$

where S_k is a suitable threshold level determined by the data of the problem. In other words, when stock falls below the threshold S_k , order just enough to bring stock up to S_k .

The preceding example illustrates the main ingredients of our formulation:

- (a) A *discrete-time system* of the form

$$x_{k+1} = f_k(x_k, u_k, w_k),$$

where f_k is some function; for example in the inventory case, we have $f_k(x_k, u_k, w_k) = x_k + u_k - w_k$.

- (b) *Independent random parameters* w_k . This will be generalized by allowing the probability distribution of w_k to depend on x_k and u_k ; in the context of the inventory example, we can think of a case where the level of demand w_k is affected by the current stock level x_k .
- (c) A *control constraint*; in the example, we have $u_k \geq 0$. In general, the constraint set will depend on x_k and the time index k , that is, $u_k \in U_k(x_k)$. To see how constraints dependent on x_k can arise in the inventory context, think of a situation where there is an upper bound B on the level of stock that can be accommodated, so $u_k \leq B - x_k$.
- (d) An *additive cost* of the form

$$E \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right\},$$

where g_k are some functions; in the inventory example, we have $g_N(x_N) = R(x_N)$ and $g_k(x_k, u_k, w_k) = r(x_k) + cu_k$.

- (e) *Optimization over (closed-loop) policies*, i.e., rules for choosing u_k for each k and each possible value of x_k .

We next consider the important special case where in addition to discrete time, the problem has a discrete state structure.

1.1.2 Discrete-State and Finite-State Problems

In the preceding example, the state x_k was a continuous real variable, and it is easy to think of multidimensional generalizations where the state is

an n -dimensional vector of real variables. It is also possible, however, that the state takes values from a discrete set, such as the integers.

A version of the inventory problem where a discrete viewpoint is more natural arises when stock is measured in whole units (such as cars), each of which is a significant fraction of x_k , u_k , or w_k . It is more appropriate then to take as state space the set of all integers rather than the set of real numbers. The form of the system equation and the cost per period will, of course, stay the same.

Generally, there are many situations where the state is naturally discrete and there is no continuous counterpart of the problem. Such situations are often conveniently specified in terms of the probabilities of transition between the states. What we need to know is $p_{ij}(u, k)$, which is the probability at time k that the next state will be j , given that the current state is i , and the control selected is u , i.e.,

$$p_{ij}(u, k) = P\{x_{k+1} = j \mid x_k = i, u_k = u\}.$$

This type of state transition can alternatively be described in terms of the discrete-time system equation

$$x_{k+1} = w_k,$$

where the probability distribution of the random parameter w_k is

$$P\{w_k = j \mid x_k = i, u_k = u\} = p_{ij}(u, k).$$

Conversely, given a discrete-state system in the form

$$x_{k+1} = f_k(x_k, u_k, w_k),$$

together with the probability distribution $P_k(w_k \mid x_k, u_k)$ of w_k , we can provide an equivalent transition probability description. The corresponding transition probabilities are given by

$$p_{ij}(u, k) = P_k\{W_k(i, u, j) \mid x_k = i, u_k = u\},$$

where $W(i, u, j)$ is the set

$$W_k(i, u, j) = \{w \mid j = f_k(i, u, w)\}.$$

Thus a discrete-state system can equivalently be described in terms of a difference equation or in terms of transition probabilities. Depending on the given problem, it may be notationally or mathematically more convenient to use one description over the other.

The following examples illustrate discrete-state problems. The first example involves a *deterministic* problem, i.e., a problem where there is no stochastic uncertainty. In such a problem, when a control is chosen at a given state, the next state is fully determined; i.e., for any state i , control u , and time k , the transition probability $p_{ij}(u, k)$ is equal to 1 for a single state j , and it is 0 for all other candidate next states. The other three examples involve stochastic problems, where the next state resulting from a given choice of control at a given state cannot be determined a priori.

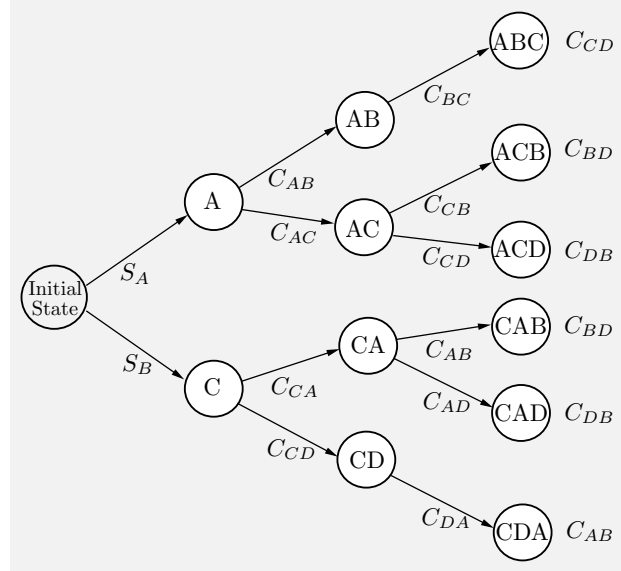


Figure 1.1.3 The transition graph of the deterministic scheduling problem of Example 1.1.2. Each arc of the graph corresponds to a decision leading from some state (the start node of the arc) to some other state (the end node of the arc). The corresponding cost is shown next to the arc. The cost of the last operation is shown as a terminal cost next to the terminal nodes of the graph.

Example 1.1.2 (A Deterministic Scheduling Problem)

Suppose that to produce a certain product, four operations must be performed on a certain machine. The operations are denoted by A, B, C, and D. We assume that operation B can be performed only after operation A has been performed, and operation D can be performed only after operation C has been performed. (Thus the sequence CDAB is allowable but the sequence CDBA is not.) The setup cost C_{mn} for passing from any operation m to any other operation n is given. There is also an initial startup cost S_A or S_C for starting with operation A or C, respectively (cf. Fig. 1.1.3). The cost of a sequence is the sum of the setup costs associated with it; for example, the operation sequence ACDB has cost

$$S_A + C_{AC} + C_{CD} + C_{DB}.$$

We can view this problem as a sequence of three decisions, namely the choice of the first three operations to be performed (the last operation is determined from the preceding three). It is appropriate to consider as state the set of operations already performed, the initial state being an artificial state corresponding to the beginning of the decision process. The possible state transitions corresponding to the possible states and decisions for this problem are shown in Fig. 1.1.3. Here the problem is deterministic, i.e., at

a given state, each choice of control leads to a uniquely determined state. For example, at state AC the decision to perform operation D leads to state ACD with certainty, and has cost C_{CD} . Deterministic problems with a finite number of states can be conveniently represented in terms of transition graphs such as the one of Fig. 1.1.3. The optimal solution corresponds to the path that starts at the initial state and ends at some state at the terminal time and has minimum sum of arc costs plus the terminal cost. We will study systematically problems of this type in Chapter 2.

Example 1.1.3 (Machine Replacement)

Consider a problem of operating efficiently over N time periods a machine that can be in any one of n states, denoted $1, 2, \dots, n$. We denote by $g(i)$ the operating cost per period when the machine is in state i , and we assume that

$$g(1) \leq g(2) \leq \dots \leq g(n).$$

The implication here is that state i is better than state $i + 1$, and state 1 corresponds to a machine in best condition.

During a period of operation, the state of the machine can become worse or it may stay unchanged. We thus assume that the transition probabilities

$$p_{ij} = P\{\text{next state will be } j \mid \text{current state is } i\}$$

satisfy

$$p_{ij} = 0 \quad \text{if } j < i.$$

We assume that at the start of each period we know the state of the machine and we must choose one of the following two options:

- (a) Let the machine operate one more period in the state it currently is.
- (b) Repair the machine and bring it to the best state 1 at a cost R .

We assume that the machine, once repaired, is guaranteed to stay in state 1 for one period. In subsequent periods, it may deteriorate to states $j > 1$ according to the transition probabilities p_{1j} .

Thus the objective here is to decide on the level of deterioration (state) at which it is worth paying the cost of machine repair, thereby obtaining the benefit of smaller future operating costs. Note that the decision should also be affected by the period we are in. For example, we would be less inclined to repair the machine when there are few periods left.

The system evolution for this problem can be described by the graphs of Fig. 1.1.4. These graphs depict the transition probabilities between various pairs of states for each value of the control and are known as *transition probability graphs* or simply *transition graphs*. Note that there is a different graph for each of the two controls.

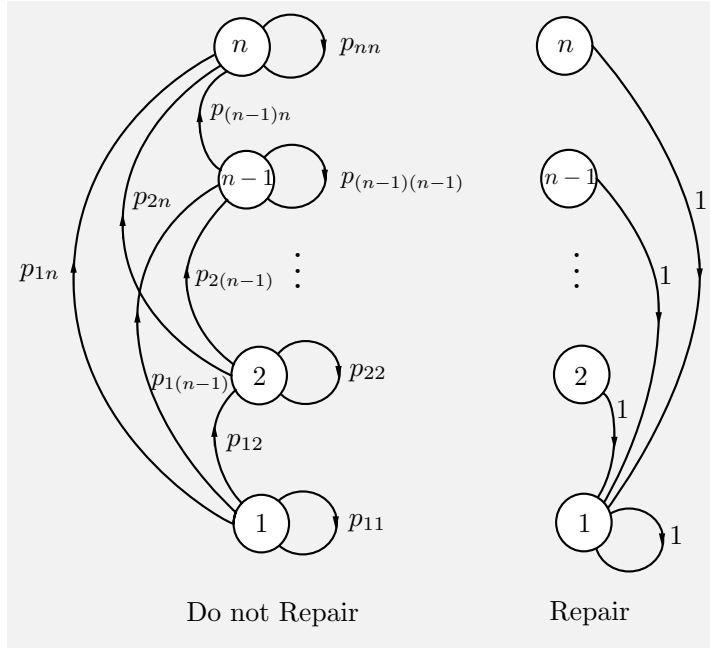


Figure 1.1.4 Machine replacement example. Transition probability graphs for each of the two possible controls (repair or not repair). At each stage and state i , the cost of repairing is $R + g(1)$, and the cost of not repairing is $g(i)$. The terminal cost is 0.

Example 1.1.4 (Control of a Queue)

Consider a queueing system with room for n customers operating over N time periods. We assume that service of a customer can start (end) only at the beginning (end) of the period and that the system can serve only one customer at a time. The probability p_m of m customer arrivals during a period is given, and the numbers of arrivals in two different periods are independent. Customers finding the system full depart without attempting to enter later. The system offers two kinds of service, *fast* and *slow*, with cost per period c_f and c_s , respectively. Service can be switched between fast and slow at the beginning of each period. With fast (slow) service, a customer in service at the beginning of a period will terminate service at the end of the period with probability q_f (respectively, q_s) independently of the number of periods the customer has been in service and the number of customers in the system ($q_f > q_s$). There is a cost $r(i)$ for each period for which there are i customers in the system. There is also a terminal cost $R(i)$ for i customers left in the system at the end of the last period.

The problem is to choose, at each period, the type of service as a function of the number of customers in the system so as to minimize the expected total cost over N periods. One expects that when there is a large number of

customers i in queue, it is better to use the fast service, and the question is to find the values of i for which this is true.

Here it is appropriate to take as state the number i of customers in the system at the start of a period and as control the type of service provided. Then, the cost per period is $r(i)$ plus c_f or c_s depending on whether fast or slow service is provided. We derive the transition probabilities of the system.

When the system is empty at the start of the period, the probability that the next state is j is independent of the type of service provided. It equals the given probability of j customer arrivals when $j < n$,

$$p_{0j}(u_f) = p_{0j}(u_s) = p_j, \quad j = 0, 1, \dots, n-1,$$

and it equals the probability of n or more customer arrivals when $j = n$,

$$p_{0n}(u_f) = p_{0n}(u_s) = \sum_{m=n}^{\infty} p_m.$$

When there is at least one customer in the system ($i > 0$), we have

$$p_{ij}(u_f) = 0, \quad \text{if } j < i-1,$$

$$p_{ij}(u_f) = q_f p_0, \quad \text{if } j = i-1,$$

$$\begin{aligned} p_{ij}(u_f) &= P\{j-i+1 \text{ arrivals, service completed}\} \\ &\quad + P\{j-i \text{ arrivals, service not completed}\} \\ &= q_f p_{j-i+1} + (1-q_f) p_{j-i}, \quad \text{if } i-1 < j < n-1, \end{aligned}$$

$$p_{i(n-1)}(u_f) = q_f \sum_{m=n-i}^{\infty} p_m + (1-q_f) p_{n-1-i},$$

$$p_{in}(u_f) = (1-q_f) \sum_{m=n-i}^{\infty} p_m.$$

The transition probabilities when slow service is provided are also given by these formulas with u_f and q_f replaced by u_s and q_s , respectively.

Example 1.1.5 (Optimizing a Chess Match Strategy)

A player is about to play a two-game chess match with an opponent, and wants to maximize his winning chances. Each game can have one of two outcomes:

- (a) A win by one of the players (1 point for the winner and 0 for the loser).
- (b) A draw (1/2 point for each of the two players).

If the score is tied at 1-1 at the end of the two games, the match goes into sudden-death mode, whereby the players continue to play until the first time

one of them wins a game (and the match). The player has two playing styles and he can choose one of the two at will in each game, independently of the style he chose in previous games.

- (1) *Timid play* with which he draws with probability $p_d > 0$, and he loses with probability $(1 - p_d)$.
- (2) *Bold play* with which he wins with probability p_w , and he loses with probability $(1 - p_w)$.

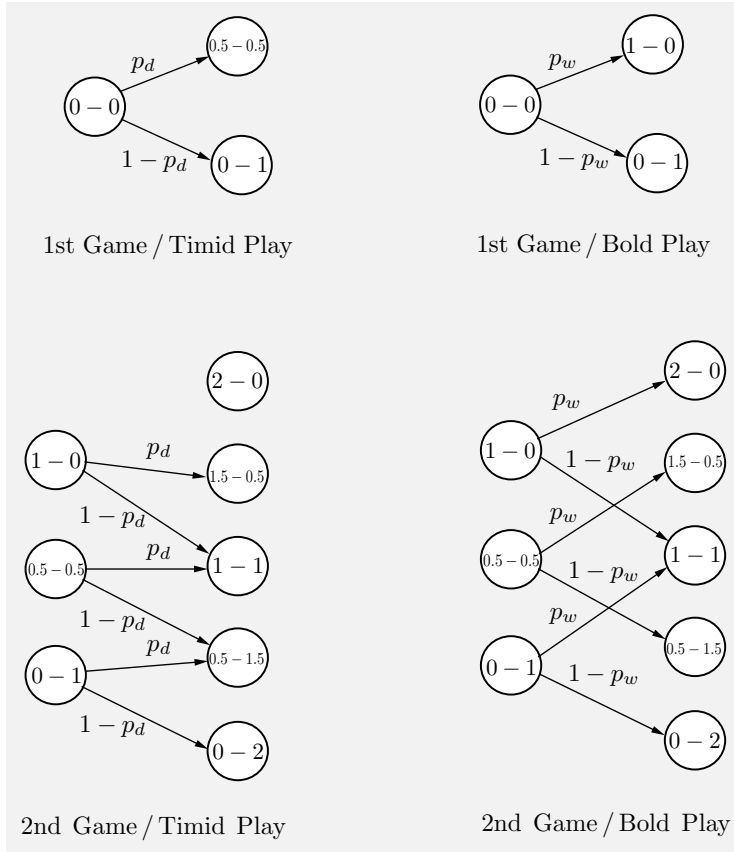


Figure 1.1.5 Chess match problem for Example 1.1.5. Transition probability graphs for each of the two possible controls (timid or bold play). Note here that the state space is not the same at each stage. The terminal cost is -1 at the winning final scores $2-0$ and $1.5-0.5$, 0 at the losing final scores $0-2$ and $0.5-1.5$, and $-p_w$ at the tied score $1-1$.

Thus, in a given game, timid play never wins, while bold play never draws. The player wants to find a style selection strategy that maximizes his proba-

bility of winning the match. Note that once the match gets into sudden death, the player should play bold, since with timid play he can at best prolong the sudden death play, while running the risk of losing. Therefore, there are only two decisions for the player to make, the selection of the playing strategy in the first two games.

We can model the problem as one with two stages, and with states the possible scores at the start of each of the first two stages (games), as shown in Fig. 1.1.5. The initial state is the initial score 0-0. The transition probabilities for each of the two different controls (playing styles) are also shown in Fig. 1.1.5. There is a cost at the terminal states: a cost of -1 at the winning scores 2-0 and 1.5-0.5, a cost of 0 at the losing scores 0-2 and 0.5-1.5, and a cost of $-p_w$ at the tied score 1-1 (since the probability of winning in sudden death is p_w). Note that to maximize the probability P of winning the match, we must minimize $-P$.

This problem has an interesting feature. One would think that if $p_w < 1/2$, the player would have a less than 50-50 chance of winning the match, even with optimal play, since his probability of losing is greater than his probability of winning any one game, regardless of his playing style. This is not so, however, because the player can adapt his playing style to the current score, but his opponent does not have that option. In other words, the player can use a closed-loop strategy, and it will be seen later that with optimal play, as determined by the DP algorithm, he has a better than 50-50 chance of winning the match provided p_w is higher than a threshold value \bar{p} , which, depending on the value of p_d , may satisfy $\bar{p} < 1/2$.

1.2 THE BASIC PROBLEM

We now formulate a general problem of decision under stochastic uncertainty over a finite number of stages. This problem, which we call *basic*, is central in this book. We will discuss solution methods based on DP in the first five chapters, and we will extend our analysis to versions of this problem involving an infinite number of stages in Chapter 5 and in Vol. II of this work.

The basic problem is very general. In particular, we will not require that the state, control, or random parameter take a finite number of values or belong to a space of n -dimensional vectors. A surprising aspect of DP is that its applicability depends very little on the nature of the state, control, and random parameter spaces. For this reason it is convenient to proceed without any assumptions on the structure of these spaces; indeed such assumptions would become a serious impediment later.

Basic Problem

We are given a discrete-time dynamic system

$$x_{k+1} = f_k(x_k, u_k, w_k), \quad k = 0, 1, \dots, N-1,$$

where the state x_k is an element of a space S_k , the control u_k is an element of a space C_k , and the random “disturbance” w_k is an element of a space D_k .

The control u_k is constrained to take values in a given nonempty subset $U(x_k) \subset C_k$, which depends on the current state x_k ; i.e., $u_k \in U_k(x_k)$ for all $x_k \in S_k$ and k .

The random disturbance w_k is characterized by a probability distribution $P_k(\cdot \mid x_k, u_k)$ that may depend explicitly on x_k and u_k but not on values of prior disturbances w_{k-1}, \dots, w_0 .

We consider the class of policies (also called control laws) that consist of a sequence of functions

$$\pi = \{\mu_0, \dots, \mu_{N-1}\},$$

where μ_k maps states x_k into controls $u_k = \mu_k(x_k)$ and is such that $\mu_k(x_k) \in U_k(x_k)$ for all $x_k \in S_k$. Such policies will be called *admissible*.

Given an initial state x_0 and an admissible policy $\pi = \{\mu_0, \dots, \mu_{N-1}\}$, the states x_k and disturbances w_k are random variables with distributions defined through the system equation

$$x_{k+1} = f_k(x_k, \mu_k(x_k), w_k), \quad k = 0, 1, \dots, N-1. \quad (1.1)$$

Thus, for given functions g_k , $k = 0, 1, \dots, N$, the expected cost of π starting at x_0 is

$$J_\pi(x_0) = E \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right\}$$

where the expectation is taken over the random variables w_k and x_k . An optimal policy π^* is one that minimizes this cost; i.e.,

$$J_{\pi^*}(x_0) = \min_{\pi \in \Pi} J_\pi(x_0),$$

where Π is the set of all admissible policies.

Note that the optimal policy π^* is associated with a fixed initial state x_0 . However, an interesting aspect of the basic problem and of DP is that it is typically possible to find a policy π^* that is simultaneously optimal for all initial states.

The optimal cost depends on x_0 and is denoted by $J^*(x_0)$; i.e.,

$$J^*(x_0) = \min_{\pi \in \Pi} J_\pi(x_0).$$

It is useful to view J^* as a function that assigns to each initial state x_0 the optimal cost $J^*(x_0)$ and call it the *optimal cost function* or *optimal value function*.[†]

The Role and Value of Information

We noted earlier the distinction between open-loop minimization, where we select all controls u_0, \dots, u_{N-1} at once at time 0, and closed-loop minimization, where we select a policy $\{\mu_0, \dots, \mu_{N-1}\}$ that applies the control $\mu_k(x_k)$ at time k with knowledge of the current state x_k (see Fig. 1.2.1). With closed-loop policies, it is possible to achieve lower cost, essentially by taking advantage of the extra information (the knowledge of the current state). The reduction in cost may be called the *value of the information* and can be significant indeed. If the information is not available, the controller cannot adapt appropriately to unexpected values of the state, and as a result the cost can be adversely affected. For example, in the inventory control example of the preceding section, the information that becomes available at the beginning of each period k is the inventory stock x_k . Clearly, this is important information to the inventory manager, who will want to adjust the amount u_k to be purchased depending on whether the current stock x_k is running high or low.

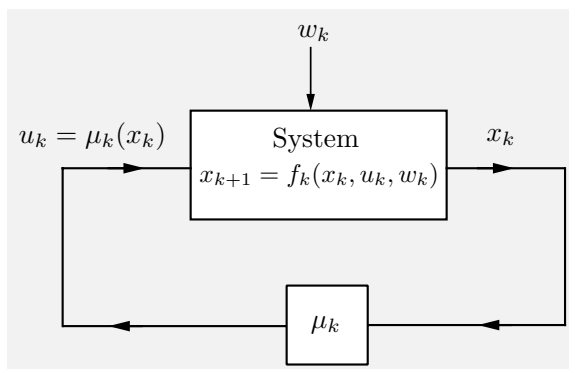


Figure 1.2.1 Information gathering in the basic problem. At each time k the controller observes the current state x_k and applies a control $u_k = \mu_k(x_k)$ that depends on that state.

[†] For the benefit of the mathematically oriented reader we note that in the preceding equation, “min” denotes the greatest lower bound (or infimum) of the set of numbers $\{J_\pi(x_0) \mid \pi \in \Pi\}$. A notation more in line with normal mathematical usage would be to write $J^*(x_0) = \inf_{\pi \in \Pi} J_\pi(x_0)$. However (as discussed in Appendix B), we find it convenient to use “min” in place of “inf” even when the infimum is not attained. It is less distracting, and it will not lead to any confusion.

Example 1.2.1

To illustrate the benefits of the proper use of information, let us consider the chess match example of the preceding section (Example 1.1.5). There, a player can select timid play (probabilities p_d and $1 - p_d$ for a draw and a loss, respectively) or bold play (probabilities p_w and $1 - p_w$ for a win and a loss, respectively) in each of the two games of the match. Suppose the player chooses a policy of playing timid if and only if he is ahead in the score, as illustrated in Fig. 1.2.2; we will see in the next section that this policy is optimal, assuming $p_d > p_w$. Then after the first game (in which he plays bold), the score is 1-0 with probability p_w and 0-1 with probability $1 - p_w$. In the second game, he plays timid in the former case and bold in the latter case. Thus after two games, the probability of a match win is $p_w p_d$, the probability of a match loss is $(1 - p_w)^2$, and the probability of a tied score is $p_w(1 - p_d) + (1 - p_w)p_w$, in which case he has a probability p_w of winning the subsequent sudden-death game. Thus the probability of winning the match with the given strategy is

$$p_w p_d + p_w(p_w(1 - p_d) + (1 - p_w)p_w),$$

which, with some rearrangement, gives

$$\text{Probability of a match win} = p_w^2(2 - p_w) + p_w(1 - p_w)p_d. \quad (1.2)$$

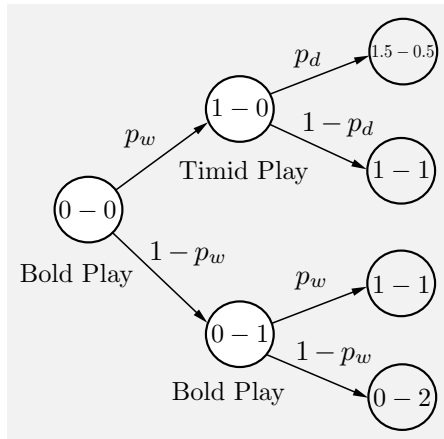


Figure 1.2.2 Illustration of the policy used in Example 1.2.1 to obtain a greater than 50-50 chance of winning the chess match and associated transition probabilities. The player chooses a policy of playing timid if and only if he is ahead in the score.

Suppose now that $p_w < 1/2$. Then the player has a greater probability of losing than winning any one game, regardless of the type of play he uses.

From this we can infer that no open-loop strategy can give the player a greater than 50-50 chance of winning the match. Yet from Eq. (1.2) it can be seen that with the closed-loop strategy of playing timid if and only if the player is ahead in the score, the chance of a match win can be greater than 50-50, provided that p_w is close enough to $1/2$ and p_d is close enough to 1. As an example, for $p_w = 0.45$ and $p_d = 0.9$, Eq. (1.2) gives a match win probability of roughly 0.53.

To calculate the value of information, let us consider the four open-loop policies, whereby we decide on the type of play to be used without waiting to see the result of the first game. These are:

- (1) Play timid in both games; this has a probability $p_d^2 p_w$ of winning the match.
- (2) Play bold in both games; this has a probability $p_w^2 + 2p_w^2(1 - p_w) = p_w^2(3 - 2p_w)$ of winning the match.
- (3) Play bold in the first game and timid in the second game; this has a probability $p_w p_d + p_w^2(1 - p_d)$ of winning the match.
- (4) Play timid in the first game and bold in the second game; this also has a probability $p_w p_d + p_w^2(1 - p_d)$ of winning the match.

The first policy is always dominated by the others, and the optimal open-loop probability of winning the match is

$$\begin{aligned} \text{Open-loop probability of win} &= \max(p_w^2(3 - 2p_w), p_w p_d + p_w^2(1 - p_d)) \\ &= p_w^2 + p_w(1 - p_w) \max(2p_w, p_d). \end{aligned} \quad (1.3)$$

Thus if $p_d > 2p_w$, we see that the optimal open-loop policy is to play timid in one of the two games and play bold in the other, while if $p_d \leq 2p_w$, it is optimal to play bold in both games.

As an example, for $p_w = 0.45$ and $p_d = 0.9$, Eq. (1.3) gives an optimal open-loop match win probability of roughly 0.425. Thus, the value of the information (the outcome of the first game) is the difference of the optimal closed-loop and open-loop values, which is approximately $0.53 - 0.425 = 0.105$. More generally, by subtracting Eqs. (1.2) and (1.3), we see that

$$\begin{aligned} \text{Value of information} &= p_w^2(2 - p_w) + p_w(1 - p_w)p_d \\ &\quad - p_w^2 - p_w(1 - p_w) \max(2p_w, p_d) \\ &= p_w(1 - p_w) \min(p_w, p_d - p_w). \end{aligned}$$

We note, however, that whereas availability of the state information cannot hurt, it may not result in an advantage either. For instance, in deterministic problems, where no random disturbances are present, one can predict the future states given the initial state and the sequence of controls. Thus, optimization over all sequences $\{u_0, u_1, \dots, u_{N-1}\}$ of controls leads to the same optimal cost as optimization over all admissible policies. The same can be true even in some stochastic problems (see for

example Exercise 1.27). This brings up a related issue. Assuming no information is forgotten, the controller actually knows the prior states and controls $x_0, u_0, \dots, x_{k-1}, u_{k-1}$ as well as the current state x_k . Therefore, the question arises whether policies that use the entire system history can be superior to policies that use just the current state. The answer turns out to be negative although the proof is technically complicated (see, e.g., [BeS78]). The intuitive reason is that, for a given time k and state x_k , all future expected costs depend explicitly just on x_k and not on prior history.

Encoding Risk in the Cost Function

As mentioned above, an important characteristic of stochastic problems is the possibility of using information with advantage. Another distinguishing characteristic is the need to take into account *risk* in the problem formulation. For example, in a typical investment problem one is not only interested in the expected profit of the investment decision, but also in its variance: given a choice between two investments with nearly equal expected profit and markedly different variance, most investors would prefer the investment with smaller variance. This indicates that expected value of cost or reward need not be the most appropriate yardstick for expressing a decision maker's preference between decisions.

As a more dramatic example of the need to take risk into account when formulating optimization problems under uncertainty, consider the so-called St. Petersburg paradox. Here, a person is offered the opportunity of paying x dollars in exchange for participation in the following game: a fair coin is flipped sequentially and the person is paid 2^k dollars, where k is the number of times heads have come up before tails come up for the first time. The decision that the person must make is whether to accept or reject participation in the game. Now if he accepts, his expected profit from the game is

$$\sum_{k=0}^{\infty} \frac{1}{2^{k+1}} \cdot 2^k - x = \infty,$$

so if his acceptance criterion is based on maximization of expected profit, he is willing to pay any amount x to enter the game. This, however, is in strong disagreement with observed behavior, due to the risk element involved in entering the game, and shows that a different formulation of the problem is needed. The formulation of problems of decision under uncertainty so that risk is properly taken into account is a deep subject with an interesting theory. An introduction to this theory is given in Appendix F. It is shown in particular that minimization of expected cost is appropriate under reasonable assumptions, provided the cost function is suitably chosen so that it properly encodes the risk preferences of the decision maker.

1.3 THE DYNAMIC PROGRAMMING ALGORITHM

The DP algorithm rests on a very simple idea, the *principle of optimality*. The name is due to Bellman, who contributed a great deal to the popularization of DP and to its transformation into a systematic tool. Roughly, the principle of optimality states the following rather obvious fact.

Principle of Optimality

Let $\pi^* = \{\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*\}$ be an optimal policy for the basic problem, and assume that when using π^* , a given state x_i occurs at time i with positive probability. Consider the subproblem whereby we are at x_i at time i and wish to minimize the “cost-to-go” from time i to time N

$$E \left\{ g_N(x_N) + \sum_{k=i}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right\}.$$

Then the truncated policy $\{\mu_i^*, \mu_{i+1}^*, \dots, \mu_{N-1}^*\}$ is optimal for this subproblem.

The intuitive justification of the principle of optimality is very simple. If the truncated policy $\{\mu_i^*, \mu_{i+1}^*, \dots, \mu_{N-1}^*\}$ were not optimal as stated, we would be able to reduce the cost further by switching to an optimal policy for the subproblem once we reach x_i . For an auto travel analogy, suppose that the fastest route from Los Angeles to Boston passes through Chicago. The principle of optimality translates to the obvious fact that the Chicago to Boston portion of the route is also the fastest route for a trip that starts from Chicago and ends in Boston.

The principle of optimality suggests that an optimal policy can be constructed in piecemeal fashion, first constructing an optimal policy for the “tail subproblem” involving the last stage, then extending the optimal policy to the “tail subproblem” involving the last two stages, and continuing in this manner until an optimal policy for the entire problem is constructed. The DP algorithm is based on this idea: it proceeds sequentially, by solving all the tail subproblems of a given time length, using the solution of the tail subproblems of shorter time length. We introduce the algorithm with two examples, one deterministic and one stochastic.

The DP Algorithm for a Deterministic Scheduling Example

Let us consider the scheduling Example 1.1.2, and let us apply the principle of optimality to calculate the optimal schedule. We have to schedule optimally the four operations A, B, C, and D. The numerical values of the transition and setup costs are shown in Fig. 1.3.1 next to the corresponding arcs.

According to the principle of optimality, the “tail” portion of an optimal schedule must be optimal. For example, suppose that the optimal schedule is CABD. Then, having scheduled first C and then A, it must be optimal to complete the schedule with BD rather than with DB. With this in mind, we solve all possible tail subproblems of length two, then all tail subproblems of length three, and finally the original problem that has length four (the subproblems of length one are of course trivial because there is only one operation that is as yet unscheduled). As we will see shortly, the tail subproblems of length $k + 1$ are easily solved once we have solved the tail subproblems of length k , and this is the essence of the DP technique.

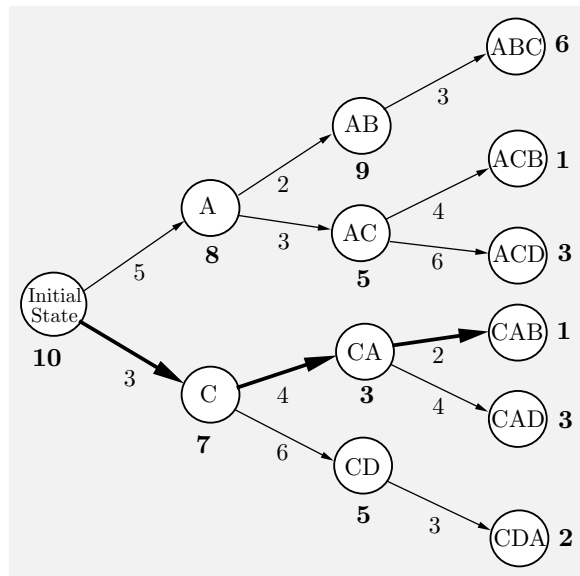


Figure 1.3.1 Transition graph of the deterministic scheduling problem, with the cost of each decision shown next to the corresponding arc. Next to each node/state we show the cost to optimally complete the schedule starting from that state. This is the optimal cost of the corresponding tail subproblem (cf. the principle of optimality). The optimal cost for the original problem is equal to 10, as shown next to the initial state. The optimal schedule corresponds to the thick-line arcs.

Tail Subproblems of Length 2: These subproblems are the ones that involve two unscheduled operations and correspond to the states AB, AC, CA, and CD (see Fig. 1.3.1)

State AB: Here it is only possible to schedule operation C as the next operation, so the optimal cost of this subproblem is 9 (the cost of

scheduling C after B, which is 3, plus the cost of scheduling D after C, which is 6).

State AC: Here the possibilities are to (a) schedule operation B and then D, which has cost 5, or (b) schedule operation D and then B, which has cost 9. The first possibility is optimal, and the corresponding cost of the tail subproblem is 5, as shown next to node AC in Fig. 1.3.1.

State CA: Here the possibilities are to (a) schedule operation B and then D, which has cost 3, or (b) schedule operation D and then B, which has cost 7. The first possibility is optimal, and the corresponding cost of the tail subproblem is 3, as shown next to node CA in Fig. 1.3.1.

State CD: Here it is only possible to schedule operation A as the next operation, so the optimal cost of this subproblem is 5.

Tail Subproblems of Length 3: These subproblems can now be solved using the optimal costs of the subproblems of length 2.

State A: Here the possibilities are to (a) schedule next operation B (cost 2) and then solve optimally the corresponding subproblem of length 2 (cost 9, as computed earlier), a total cost of 11, or (b) schedule next operation C (cost 3) and then solve optimally the corresponding subproblem of length 2 (cost 5, as computed earlier), a total cost of 8. The second possibility is optimal, and the corresponding cost of the tail subproblem is 8, as shown next to node A in Fig. 1.3.1.

State C: Here the possibilities are to (a) schedule next operation A (cost 4) and then solve optimally the corresponding subproblem of length 2 (cost 3, as computed earlier), a total cost of 7, or (b) schedule next operation D (cost 6) and then solve optimally the corresponding subproblem of length 2 (cost 5, as computed earlier), a total cost of 11. The first possibility is optimal, and the corresponding cost of the tail subproblem is 7, as shown next to node C in Fig. 1.3.1.

Original Problem of Length 4: The possibilities here are (a) start with operation A (cost 5) and then solve optimally the corresponding subproblem of length 3 (cost 8, as computed earlier), a total cost of 13, or (b) start with operation C (cost 3) and then solve optimally the corresponding subproblem of length 3 (cost 7, as computed earlier), a total cost of 10. The second possibility is optimal, and the corresponding optimal cost is 10, as shown next to the initial state node in Fig. 1.3.1.

Note that having computed the optimal cost of the original problem through the solution of all the tail subproblems, we can construct the optimal schedule by starting at the initial node and proceeding forward, each time choosing the operation that starts the optimal schedule for the cor-

responding tail subproblem. In this way, by inspection of the graph and the computational results of Fig. 1.3.1, we determine that CABD is the optimal schedule.

The DP Algorithm for the Inventory Control Example

Consider the inventory control Example 1.1.1. Similar to the solution of the preceding deterministic scheduling problem, we calculate sequentially the optimal costs of all the tail subproblems, going from shorter to longer problems. The only difference is that the optimal costs are computed as expected values, since the problem here is stochastic.

Tail Subproblems of Length 1: Assume that at the beginning of period $N - 1$ the stock is x_{N-1} . Clearly, no matter what happened in the past, the inventory manager should order the amount of inventory that minimizes over $u_{N-1} \geq 0$ the sum of the ordering cost and the expected terminal holding/shortage cost. Thus, he should minimize over u_{N-1} the sum $cu_{N-1} + E\{R(x_N)\}$, which can be written as

$$cu_{N-1} + \min_{w_{N-1}} E \{R(x_{N-1} + u_{N-1} - w_{N-1})\}.$$

Adding the holding/shortage cost of period $N - 1$, we see that the optimal cost for the last period (plus the terminal cost) is given by

$$J_{N-1}(x_{N-1}) = r(x_{N-1}) + \min_{u_{N-1} \geq 0} \left[cu_{N-1} + \min_{w_{N-1}} E \{R(x_{N-1} + u_{N-1} - w_{N-1})\} \right].$$

Naturally, J_{N-1} is a function of the stock x_{N-1} . It is calculated either analytically or numerically (in which case a table is used for computer storage of the function J_{N-1}). In the process of calculating J_{N-1} , we obtain the optimal inventory policy $\mu_{N-1}^*(x_{N-1})$ for the last period: $\mu_{N-1}^*(x_{N-1})$ is the value of u_{N-1} that minimizes the right-hand side of the preceding equation for a given value of x_{N-1} .

Tail Subproblems of Length 2: Assume that at the beginning of period $N - 2$ the stock is x_{N-2} . It is clear that the inventory manager should order the amount of inventory that minimizes not just the expected cost of period $N - 2$ but rather the

$$\begin{aligned} & (\text{expected cost of period } N - 2) + (\text{expected cost of period } N - 1, \\ & \text{given that an optimal policy will be used at period } N - 1), \end{aligned}$$

which is equal to

$$r(x_{N-2}) + cu_{N-2} + E\{J_{N-1}(x_{N-1})\}.$$

Using the system equation $x_{N-1} = x_{N-2} + u_{N-2} - w_{N-2}$, the last term is also written as $J_{N-1}(x_{N-2} + u_{N-2} - w_{N-2})$.

Thus the optimal cost for the last two periods given that we are at state x_{N-2} , denoted $J_{N-2}(x_{N-2})$, is given by

$$J_{N-2}(x_{N-2}) = r(x_{N-2}) + \min_{u_{N-2} \geq 0} \left[cu_{N-2} + E_{w_{N-2}} \{ J_{N-1}(x_{N-2} + u_{N-2} - w_{N-2}) \} \right]$$

Again $J_{N-2}(x_{N-2})$ is calculated for every x_{N-2} . At the same time, the optimal policy $\mu_{N-2}^*(x_{N-2})$ is also computed.

Tail Subproblems of Length $N - k$: Similarly, we have that at period k , when the stock is x_k , the inventory manager should order u_k to minimize

(expected cost of period k) + (expected cost of periods $k + 1, \dots, N - 1$,
given that an optimal policy will be used for these periods).

By denoting by $J_k(x_k)$ the optimal cost, we have

$$J_k(x_k) = r(x_k) + \min_{u_k \geq 0} \left[cu_k + E_{w_k} \{ J_{k+1}(x_k + u_k - w_k) \} \right], \quad (1.4)$$

which is actually the DP equation for this problem.

The functions $J_k(x_k)$ denote the optimal expected cost for the tail subproblem that starts at period k with initial inventory x_k . These functions are computed recursively backward in time, starting at period $N - 1$ and ending at period 0. The value $J_0(x_0)$ is the optimal expected cost when the initial stock at time 0 is x_0 . During the calculations, the optimal policy is simultaneously computed from the minimization in the right-hand side of Eq. (1.4).

The example illustrates the main advantage offered by DP. While the original inventory problem requires an optimization over the set of policies, the DP algorithm of Eq. (1.4) decomposes this problem into a sequence of minimizations carried out over the set of controls. Each of these minimizations is much simpler than the original problem.

The DP Algorithm

We now state the DP algorithm for the basic problem and show its optimality by translating into mathematical terms the heuristic argument given above for the inventory example.

Proposition 1.3.1: For every initial state x_0 , the optimal cost $J^*(x_0)$ of the basic problem is equal to $J_0(x_0)$, given by the last step of the following algorithm, which proceeds backward in time from period $N - 1$ to period 0:

$$J_N(x_N) = g_N(x_N), \quad (1.5)$$

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right\}, \quad k = 0, 1, \dots, N - 1, \quad (1.6)$$

where the expectation is taken with respect to the probability distribution of w_k , which depends on x_k and u_k . Furthermore, if $u_k^* = \mu_k^*(x_k)$ minimizes the right side of Eq. (1.6) for each x_k and k , the policy $\pi^* = \{\mu_0^*, \dots, \mu_{N-1}^*\}$ is optimal.

Proof:[†] For any admissible policy $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$ and each $k = 0, 1, \dots, N - 1$, denote $\pi^k = \{\mu_k, \mu_{k+1}, \dots, \mu_{N-1}\}$. For $k = 0, 1, \dots, N - 1$, let $J_k^*(x_k)$ be the optimal cost for the $(N - k)$ -stage problem that starts at state x_k and time k , and ends at time N ,

$$J_k^*(x_k) = \min_{\pi^k} E_{w_k, \dots, w_{N-1}} \left\{ g_N(x_N) + \sum_{i=k}^{N-1} g_i(x_i, \mu_i(x_i), w_i) \right\}.$$

For $k = N$, we define $J_N^*(x_N) = g_N(x_N)$. We will show by induction that the functions J_k^* are equal to the functions J_k generated by the DP algorithm, so that for $k = 0$, we will obtain the desired result.

Indeed, we have by definition $J_N^* = J_N = g_N$. Assume that for some k and all x_{k+1} , we have $J_{k+1}^*(x_{k+1}) = J_{k+1}(x_{k+1})$. Then, since $\pi^k = (\mu_k, \pi^{k+1})$, we have for all x_k

$$J_k^*(x_k) = \min_{(\mu_k, \pi^{k+1})} E_{w_k, \dots, w_{N-1}} \left\{ g_k(x_k, \mu_k(x_k), w_k) + g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), w_i) \right\}$$

[†] Our proof is somewhat informal and assumes that the functions J_k are well-defined and finite. For a strictly rigorous proof, some technical mathematical issues must be addressed; see Section 1.5. These issues do not arise if the disturbance w_k takes a finite or countable number of values and the expected values of all terms in the expression of the cost function (1.1) are well-defined and finite for every admissible policy π .

$$\begin{aligned}
&= \min_{\mu_k} E_{w_k} \left\{ g_k(x_k, \mu_k(x_k), w_k) \right. \\
&\quad \left. + \min_{\pi^{k+1}} E_{w_{k+1}, \dots, w_{N-1}} \left\{ g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), w_i) \right\} \right\} \\
&= \min_{\mu_k} E_{w_k} \left\{ g_k(x_k, \mu_k(x_k), w_k) + J_{k+1}^*(f_k(x_k, \mu_k(x_k), w_k)) \right\} \\
&= \min_{\mu_k} E_{w_k} \left\{ g_k(x_k, \mu_k(x_k), w_k) + J_{k+1}(f_k(x_k, \mu_k(x_k), w_k)) \right\} \\
&= \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right\} \\
&= J_k(x_k),
\end{aligned}$$

completing the induction. In the second equation above, we moved the minimum over π^{k+1} inside the braced expression, using a principle of optimality argument: “the tail portion of an optimal policy is optimal for the tail subproblem” (a more rigorous justification of this step is given in Section 1.5). In the third equation, we used the definition of J_{k+1}^* , and in the fourth equation we used the induction hypothesis. In the fifth equation, we converted the minimization over μ_k to a minimization over u_k , using the fact that for any function F of x and u , we have

$$\min_{\mu \in M} F(x, \mu(x)) = \min_{u \in U(x)} F(x, u),$$

where M is the set of all functions $\mu(x)$ such that $\mu(x) \in U(x)$ for all x . **Q.E.D.**

The argument of the preceding proof provides an interpretation of $J_k(x_k)$ as the optimal cost for an $(N - k)$ -stage problem starting at state x_k and time k , and ending at time N . We consequently call $J_k(x_k)$ the *cost-to-go* at state x_k and time k , and refer to J_k as the *cost-to-go function* or *optimal cost function* at time k .[†]

Ideally, we would like to use the DP algorithm to obtain closed-form expressions for J_k or an optimal policy. In this book, we will discuss quite a few models that admit analytical solution by DP. Even if such models rely on oversimplified assumptions, they are often very useful. They may provide valuable insights about the structure of the optimal solution of more complex models, and they may form the basis for suboptimal control schemes. Furthermore, the broad collection of analytically solvable models provides helpful guidelines for modeling: when faced with a new problem it

[†] In maximization problems the DP algorithm (1.6) is written with maximization in place of minimization, and then J_k is referred to as the *optimal value function* at time k .

is worth trying to pattern its model after one of the principal analytically tractable models.

Unfortunately, in many practical cases an analytical solution is not possible, and one has to resort to numerical execution of the DP algorithm. This may be quite time-consuming since the minimization in the DP Eq. (1.6) must be carried out for each value of x_k . The state space must be discretized in some way if it is not already a finite set. The computational requirements are proportional to the number of possible values of x_k , so for complex problems the computational burden may be excessive. Nonetheless, DP is the only general approach for sequential optimization under uncertainty, and even when it is computationally prohibitive, it can serve as the basis for more practical suboptimal approaches, which will be discussed in Chapter 6. Moreover, the DP computation is still far more economical than a brute force search, as the following example illustrates.

Example 1.3.1 (Complexity Aspects of DP)

Let us calculate more precisely the computational requirements of DP in a finite-spaces context. Assume that the state spaces X_0, X_1, \dots, X_{N-1} have no more than n elements each, and that at each state there are no more than m control elements available. Then the total number of state-control-time triples is no more than nmN . Thus nmN is an upper bound to the total number of times that expressions of the form

$$E_{w_k} \left\{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right\},$$

need to be calculated in the course of the DP algorithm [cf. Eq. (1.6)].

Of course the preceding expression may involve potentially significant computation. In particular, the expected value requires a number of calculations of the form

$$g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)). \quad (1.7)$$

This number is between 1 (for a deterministic problem) to n (typically, for a stochastic problem if the distribution of w_k is known). In either case, the number of times the expression (1.7) needs to be calculated is polynomial in n , m , and N .

Let us compare this computation with a brute force approach, which enumerates and compares all the possible solutions. A closed-loop policy $\{\mu_0, \dots, \mu_{N-1}\}$ is characterized by a single control at each state-time pair (x_k, k) , and there are no more than nN such pairs. With as many as m controls available at each of these pairs, we see that the number of distinct policies can be as many as m^{nN} . If we restrict ourselves to open-loop sequences (which we can for deterministic problems), still for a given initial condition, the number of possible sequences is as many as m^N . Thus the size of the solution space grows exponentially with N for both stochastic and deterministic problems.

A final observation is that the favorable complexity properties of DP depend critically on the additive structure of the cost function. We will see later in Section 1.4 that we can convert problems with a nonadditive cost structure to the basic problem format through a technique called *state augmentation*. However, in doing so the number of states grows again exponentially with N .

Dynamic Programming Examples

Let us now illustrate some of the analytical and computational aspects of DP by means of examples.

Example 1.3.2

We will go through the details of the DP algorithm for a stochastic inventory control problem that is similar to the one of Sections 1.1 and 1.2, but slightly different in some details. In particular, we assume that the inventory u_k and the demand w_k are nonnegative integers, and that the excess demand $(w_k - x_k - u_k)$ is lost. As a result, the stock equation takes the form

$$x_{k+1} = \max(0, x_k + u_k - w_k).$$

We also assume that there is an upper bound of 2 units on the stock that can be stored, i.e. there is a constraint $x_k + u_k \leq 2$. The holding/storage cost for the k th period is given by

$$(x_k + u_k - w_k)^2,$$

implying a penalty both for excess inventory and for unmet demand at the end of the k th period. The ordering cost is 1 per unit stock ordered. Thus the cost per period is

$$g_k(x_k, u_k, w_k) = u_k + (x_k + u_k - w_k)^2.$$

The terminal cost is assumed to be 0,

$$g_N(x_N) = 0.$$

The planning horizon N is 3 periods, and the initial stock x_0 is 0. The demand w_k has the same probability distribution for all periods, given by

$$p(w_k = 0) = 0.1, \quad p(w_k = 1) = 0.7, \quad p(w_k = 2) = 0.2.$$

The system can also be represented in terms of the transition probabilities $p_{ij}(u)$ between the three possible states, for the different values of the control (see Fig. 1.3.2).

The starting equation for the DP algorithm is

$$J_3(x_3) = 0,$$

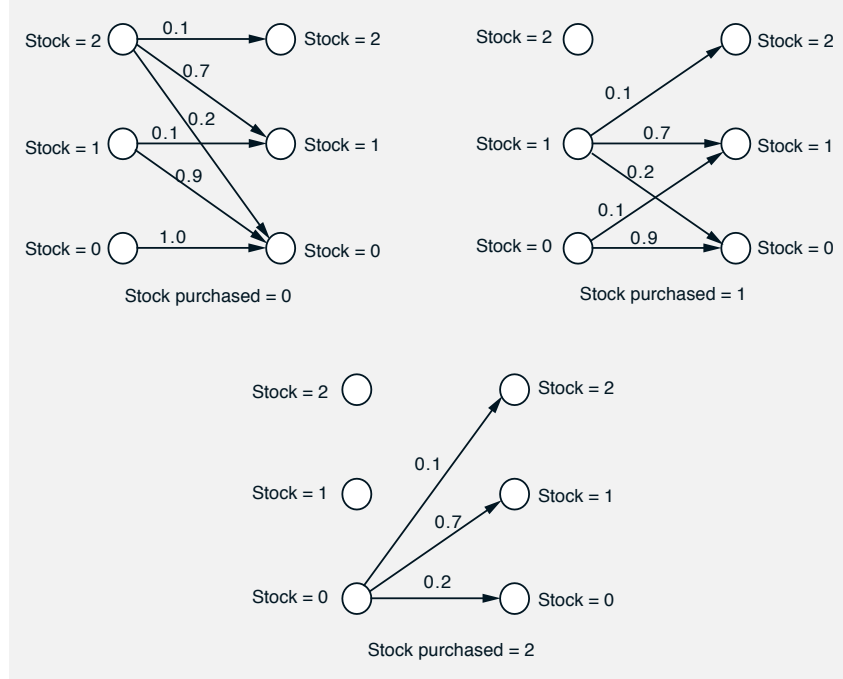


Figure 1.3.2 System and DP results for Example 1.3.2. The transition probability diagrams for the different values of stock purchased (control) are shown. The numbers next to the arcs are the transition probabilities. The control $u = 1$ is not available at state 2 because of the limitation $x_k + u_k \leq 2$. Similarly, the control $u = 2$ is not available at states 1 and 2. The results of the DP algorithm are given in the table.

since the terminal state cost is 0 [cf. Eq. (1.5)]. The algorithm takes the form [cf. Eq. (1.6)]

$$J_k(x_k) = \min_{\substack{0 \leq u_k \leq 2 - x_k \\ u_k = 0, 1, 2}} E_{w_k} \left\{ u_k + (x_k + u_k - w_k)^2 + J_{k+1}(\max(0, x_k + u_k - w_k)) \right\},$$

where $k = 0, 1, 2$, and x_k, u_k, w_k can take the values 0, 1, and 2.

Period 2: We compute $J_2(x_2)$ for each of the three possible states. We have

$$\begin{aligned} J_2(0) &= \min_{u_2=0,1,2} E \left\{ u_2 + (u_2 - w_2)^2 \right\} \\ &= \min_{u_2=0,1,2} \left[u_2 + 0.1(u_2)^2 + 0.7(u_2 - 1)^2 + 0.2(u_2 - 2)^2 \right]. \end{aligned}$$

We calculate the expectation of the right side for each of the three possible values of u_2 :

$$\begin{aligned} u_2 = 0 : E\{\cdot\} &= 0.7 \cdot 1 + 0.2 \cdot 4 = 1.5, \\ u_2 = 1 : E\{\cdot\} &= 1 + 0.1 \cdot 1 + 0.2 \cdot 1 = 1.3, \\ u_2 = 2 : E\{\cdot\} &= 2 + 0.1 \cdot 4 + 0.7 \cdot 1 = 3.1. \end{aligned}$$

Hence we have, by selecting the minimizing u_2 ,

$$J_2(0) = 1.3, \quad \mu_2^*(0) = 1.$$

For $x_2 = 1$, we have

$$\begin{aligned} J_2(1) &= \min_{u_2=0,1} E \left\{ u_2 + (1 + u_2 - w_2)^2 \right\} \\ &= \min_{u_2=0,1} \left[u_2 + 0.1(1 + u_2)^2 + 0.7(u_2)^2 + 0.2(u_2 - 1)^2 \right]. \end{aligned}$$

The expected value in the right side is

$$\begin{aligned} u_2 = 0 : E\{\cdot\} &= 0.1 \cdot 1 + 0.2 \cdot 1 = 0.3, \\ u_2 = 1 : E\{\cdot\} &= 1 + 0.1 \cdot 4 + 0.7 \cdot 1 = 2.1. \end{aligned}$$

Hence

$$J_2(1) = 0.3, \quad \mu_2^*(1) = 0.$$

For $x_2 = 2$, the only admissible control is $u_2 = 0$, so we have

$$\begin{aligned} J_2(2) &= E \left\{ (2 - w_2)^2 \right\} = 0.1 \cdot 4 + 0.7 \cdot 1 = 1.1, \\ J_2(2) &= 1.1, \quad \mu_2^*(2) = 0. \end{aligned}$$

Period 1: Again we compute $J_1(x_1)$ for each of the three possible states $x_1 = 0, 1, 2$, using the values $J_2(0)$, $J_2(1)$, $J_2(2)$ obtained in the previous period. For $x_1 = 0$, we have

$$J_1(0) = \min_{u_1=0,1,2} E \left\{ u_1 + (u_1 - w_1)^2 + J_2(\max(0, u_1 - w_1)) \right\},$$

$$\begin{aligned} u_1 = 0 : E\{\cdot\} &= 0.1 \cdot J_2(0) + 0.7(1 + J_2(0)) + 0.2(4 + J_2(0)) = 2.8, \\ u_1 = 1 : E\{\cdot\} &= 1 + 0.1(1 + J_2(1)) + 0.7 \cdot J_2(0) + 0.2(1 + J_2(0)) = 2.5, \\ u_1 = 2 : E\{\cdot\} &= 2 + 0.1(4 + J_2(2)) + 0.7(1 + J_2(1)) + 0.2 \cdot J_2(0) = 3.68, \end{aligned}$$

$$J_1(0) = 2.5, \quad \mu_1^*(0) = 1.$$

For $x_1 = 1$, we have

$$J_1(1) = \min_{u_1=0,1} E \left\{ u_1 + (1 + u_1 - w_1)^2 + J_2(\max(0, 1 + u_1 - w_1)) \right\},$$

$$u_1 = 0 : E\{\cdot\} = 0.1(1 + J_2(1)) + 0.7 \cdot J_2(0) + 0.2(1 + J_2(0)) = 1.5,$$

$$u_1 = 1 : E\{\cdot\} = 1 + 0.1(4 + J_2(2)) + 0.7(1 + J_2(1)) + 0.2 \cdot J_2(0) = 2.68,$$

$$J_1(1) = 1.5, \quad \mu_1^*(1) = 0.$$

For $x_1 = 2$, the only admissible control is $u_1 = 0$, so we have

$$\begin{aligned} J_1(2) &= E_{w_1} \left\{ (2 - w_1)^2 + J_2(\max(0, 2 - w_1)) \right\} \\ &= 0.1(4 + J_2(2)) + 0.7(1 + J_2(1)) + 0.2 \cdot J_2(0) \\ &= 1.68, \end{aligned}$$

$$J_1(2) = 1.68, \quad \mu_1^*(2) = 0.$$

Period 0: Here we need to compute only $J_0(0)$ since the initial state is known to be 0. We have

$$J_0(0) = \min_{u_0=0,1,2} E \left\{ u_0 + (u_0 - w_0)^2 + J_1(\max(0, u_0 - w_0)) \right\},$$

$$u_0 = 0 : E\{\cdot\} = 0.1 \cdot J_1(0) + 0.7(1 + J_1(0)) + 0.2(4 + J_1(0)) = 4.0,$$

$$u_0 = 1 : E\{\cdot\} = 1 + 0.1(1 + J_1(1)) + 0.7 \cdot J_1(0) + 0.2(1 + J_1(0)) = 3.7,$$

$$u_0 = 2 : E\{\cdot\} = 2 + 0.1(4 + J_1(2)) + 0.7(1 + J_1(1)) + 0.2 \cdot J_1(0) = 4.818,$$

$$J_0(0) = 3.7, \quad \mu_0^*(0) = 1.$$

If the initial state were not known a priori, we would have to compute in a similar manner $J_0(1)$ and $J_0(2)$, as well as the minimizing u_0 . The reader may verify (Exercise 1.1) that these calculations yield

$$J_0(1) = 2.7, \quad \mu_0^*(1) = 0,$$

$$J_0(2) = 2.818, \quad \mu_0^*(2) = 0.$$

Thus the optimal ordering policy for each period is to order one unit if the current stock is zero and order nothing otherwise. The results of the DP algorithm are given in tabular form in Fig. 1.3.2.

Example 1.3.3 (A Linear-Quadratic Problem)

This is an example involving a one-dimensional linear system and a quadratic cost function. It illustrates an important class of problems that admit an analytical solution, and will be discussed in much greater detail later.

A certain material is passed through a sequence of two ovens (see Fig. 1.3.3). Denote

x_0 : initial temperature of the material,

$x_k, k = 1, 2$: temperature of the material at the exit of oven k ,

$u_{k-1}, k = 1, 2$: prevailing temperature in oven k .

We assume a model of the form

$$x_{k+1} = (1 - a)x_k + au_k, \quad k = 0, 1,$$

where a is a known scalar from the interval $(0, 1)$. The objective is to get the final temperature x_2 close to a given target T , while expending relatively little energy. This is expressed by a cost function of the form

$$r(x_2 - T)^2 + u_0^2 + u_1^2,$$

where $r > 0$ is a given scalar. We assume no constraints on u_k . (In reality, there are constraints, but if we can solve the unconstrained problem and verify that the solution satisfies the constraints, everything will be fine.) The problem is deterministic; i.e., there is no stochastic uncertainty. However, such problems can be placed within the basic framework by introducing a fictitious disturbance taking a unique value with probability one.

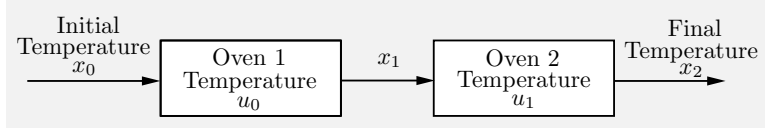


Figure 1.3.3 The linear-quadratic problem of Example 1.3.3. The temperature of the material evolves according to $x_{k+1} = (1 - a)x_k + au_k$, where a is some scalar with $0 < a < 1$.

We have $N = 2$ and a terminal cost $g_2(x_2) = r(x_2 - T)^2$, so the initial condition for the DP algorithm is [cf. Eq. (1.5)]

$$J_2(x_2) = r(x_2 - T)^2.$$

For the next-to-last stage, we have [cf. Eq. (1.6)]

$$\begin{aligned} J_1(x_1) &= \min_{u_1} [u_1^2 + J_2(x_2)] \\ &= \min_{u_1} \left[u_1^2 + J_2((1 - a)x_1 + au_1) \right]. \end{aligned}$$

Substituting the previous form of J_2 , we obtain

$$J_1(x_1) = \min_{u_1} \left[u_1^2 + r((1-a)x_1 + au_1 - T)^2 \right]. \quad (1.8)$$

This minimization will be done by setting to zero the derivative with respect to u_1 . This yields

$$0 = 2u_1 + 2ra((1-a)x_1 + au_1 - T),$$

and by collecting terms and solving for u_1 , we obtain the optimal temperature for the last oven:

$$\mu_1^*(x_1) = \frac{ra(T - (1-a)x_1)}{1 + ra^2}.$$

Note that this is not a single control but rather a control function, a rule that tells us the optimal oven temperature $u_1 = \mu_1^*(x_1)$ for each possible value of the state x_1 .

By substituting the optimal u_1 in the expression (1.8) for J_1 , we obtain

$$\begin{aligned} J_1(x_1) &= \frac{r^2 a^2 ((1-a)x_1 - T)^2}{(1 + ra^2)^2} + r \left((1-a)x_1 + \frac{ra^2(T - (1-a)x_1)}{1 + ra^2} - T \right)^2 \\ &= \frac{r^2 a^2 ((1-a)x_1 - T)^2}{(1 + ra^2)^2} + r \left(\frac{ra^2}{1 + ra^2} - 1 \right)^2 ((1-a)x_1 - T)^2 \\ &= \frac{r((1-a)x_1 - T)^2}{1 + ra^2}. \end{aligned}$$

We now go back one stage. We have [cf. Eq. (1.6)]

$$J_0(x_0) = \min_{u_0} [u_0^2 + J_1(x_1)] = \min_{u_0} \left[u_0^2 + J_1((1-a)x_0 + au_0) \right],$$

and by substituting the expression already obtained for J_1 , we have

$$J_0(x_0) = \min_{u_0} \left[u_0^2 + \frac{r((1-a)^2 x_0 + (1-a)au_0 - T)^2}{1 + ra^2} \right].$$

We minimize with respect to u_0 by setting the corresponding derivative to zero. We obtain

$$0 = 2u_0 + \frac{2r(1-a)a((1-a)^2 x_0 + (1-a)au_0 - T)}{1 + ra^2}.$$

This yields, after some calculation, the optimal temperature of the first oven:

$$\mu_0^*(x_0) = \frac{r(1-a)a(T - (1-a)^2 x_0)}{1 + ra^2(1 + (1-a)^2)}.$$

The optimal cost is obtained by substituting this expression in the formula for J_0 . This leads to a straightforward but lengthy calculation, which in the end yields the rather simple formula

$$J_0(x_0) = \frac{r((1-a)^2x_0 - T)^2}{1 + ra^2(1 + (1-a)^2)}.$$

This completes the solution of the problem.

One noteworthy feature in this example is the facility with which we obtained an analytical solution. A little thought while tracing the steps of the algorithm will convince the reader that what simplifies the solution is the quadratic nature of the cost and the linearity of the system equation. In Section 3.1 we will see that, generally, when the system is linear and the cost is quadratic, the optimal policy and cost-to-go function are given by closed-form expressions, regardless of the number of stages N .

Another noteworthy feature of the example is that the optimal policy remains unaffected when a zero-mean stochastic disturbance is added in the system equation. To see this, assume that the material's temperature evolves according to

$$x_{k+1} = (1-a)x_k + au_k + w_k, \quad k = 0, 1,$$

where w_0, w_1 are independent random variables with given distribution, zero mean

$$E\{w_0\} = E\{w_1\} = 0,$$

and finite variance. Then the equation for J_1 [cf. Eq. (1.6)] becomes

$$\begin{aligned} J_1(x_1) &= \min_{u_1} E \left\{ u_1^2 + r((1-a)x_1 + au_1 + w_1 - T)^2 \right\} \\ &= \min_{u_1} \left[u_1^2 + r((1-a)x_1 + au_1 - T)^2 \right. \\ &\quad \left. + 2rE\{w_1\}((1-a)x_1 + au_1 - T) + rE\{w_1^2\} \right]. \end{aligned}$$

Since $E\{w_1\} = 0$, we obtain

$$J_1(x_1) = \min_{u_1} \left[u_1^2 + r((1-a)x_1 + au_1 - T)^2 \right] + rE\{w_1^2\}.$$

Comparing this equation with Eq. (1.8), we see that the presence of w_1 has resulted in an additional inconsequential constant term, $rE\{w_1^2\}$. Therefore, the optimal policy for the last stage remains unaffected by the presence of w_1 , while $J_1(x_1)$ is increased by $rE\{w_1^2\}$. It can be seen that a similar situation also holds for the first stage. In particular, the optimal cost is given by the same expression as before except for an additive constant that depends on $E\{w_0^2\}$ and $E\{w_1^2\}$.

If the optimal policy is unaffected when the disturbances are replaced by their means, we say that *certainty equivalence* holds. We will derive certainty equivalence results for several types of problems involving a linear system and a quadratic cost (see Sections 3.1, 4.2, and 4.3).

Example 1.3.4 (Optimizing a Chess Match Strategy)

Consider the chess match Example 1.1.5. There, a player can select timid play (probabilities p_d and $1 - p_d$ for a draw or loss, respectively) or bold play (probabilities p_w and $1 - p_w$ for a win or loss, respectively) in each game of the match. We want to formulate a DP algorithm for finding the policy that maximizes the player's probability of winning the match. Note that here we are dealing with a maximization problem. We can convert the problem to a minimization problem by changing the sign of the cost function, but a simpler alternative, which we will generally adopt, is to replace the minimization in the DP algorithm with maximization.

Let us consider the general case of an N -game match, and let the state be the *net score*, i.e., the difference between the points of the player minus the points of the opponent (so a state of 0 corresponds to an even score). The optimal cost-to-go function at the start of the k th game is given by the DP recursion

$$J_k(x_k) = \max \left[p_d J_{k+1}(x_k) + (1 - p_d) J_{k+1}(x_k - 1), \right. \\ \left. p_w J_{k+1}(x_k + 1) + (1 - p_w) J_{k+1}(x_k - 1) \right]. \quad (1.9)$$

The maximum above is taken over the two possible decisions:

- (a) Timid play, which keeps the score at x_k with probability p_d , and changes x_k to $x_k - 1$ with probability $1 - p_d$.
- (b) Bold play, which changes x_k to $x_k + 1$ or to $x_k - 1$ with probabilities p_w or $(1 - p_w)$, respectively.

It is optimal to play bold when

$$p_w J_{k+1}(x_k + 1) + (1 - p_w) J_{k+1}(x_k - 1) \geq p_d J_{k+1}(x_k) + (1 - p_d) J_{k+1}(x_k - 1)$$

or equivalently, if

$$\frac{p_w}{p_d} \geq \frac{J_{k+1}(x_k) - J_{k+1}(x_k - 1)}{J_{k+1}(x_k + 1) - J_{k+1}(x_k - 1)}. \quad (1.10)$$

The DP recursion is started with

$$J_N(x_N) = \begin{cases} 1 & \text{if } x_N > 0, \\ p_w & \text{if } x_N = 0, \\ 0 & \text{if } x_N < 0. \end{cases} \quad (1.11)$$

In this equation, we have $J_N(0) = p_w$ because when the score is even after N games ($x_N = 0$), it is optimal to play bold in the first game of sudden death.

By executing the DP algorithm (1.9) starting with the terminal condition (1.11), and using the criterion (1.10) for optimality of bold play, we find the following, assuming that $p_d > p_w$:

$$\begin{aligned} J_{N-1}(x_{N-1}) &= 1 \text{ for } x_{N-1} > 1; && \text{optimal play: either} \\ J_{N-1}(1) &= \max[p_d + (1 - p_d)p_w, p_w + (1 - p_w)p_w] \\ &= p_d + (1 - p_d)p_w; && \text{optimal play: timid} \\ J_{N-1}(0) &= p_w; && \text{optimal play: bold} \\ J_{N-1}(-1) &= p_w^2; && \text{optimal play: bold} \\ J_{N-1}(x_{N-1}) &= 0 \text{ for } x_{N-1} < -1; && \text{optimal play: either.} \end{aligned}$$

Also, given $J_{N-1}(x_{N-1})$, and using Eqs. (1.9) and (1.10), we obtain

$$\begin{aligned} J_{N-2}(0) &= \max \left[p_d p_w + (1 - p_d) p_w^2, p_w (p_d + (1 - p_d) p_w) + (1 - p_w) p_w^2 \right] \\ &= p_w (p_w + (p_w + p_d)(1 - p_w)), \end{aligned}$$

and that if the score is even with 2 games remaining, it is optimal to play bold. Thus for a 2-game match, the optimal policy for both periods is to play timid if and only if the player is ahead in the score. The region of pairs (p_w, p_d) for which the player has a better than 50-50 chance to win a 2-game match is

$$R_2 = \left\{ (p_w, p_d) \mid J_0(0) = p_w (p_w + (p_w + p_d)(1 - p_w)) > 1/2 \right\},$$

and, as noted in Example 1.2.1, it includes points where $p_w < 1/2$.

Example 1.3.5 (Finite State Systems)

We mentioned earlier (cf. the examples in Section 1.1) that systems with a finite number of states can be represented either in terms of a discrete-time system equation or in terms of the probabilities of transition between the states. Let us work out the DP algorithm corresponding to the latter case. We assume for the sake of the following discussion that the problem is stationary, i.e., the transition probabilities, the cost per stage, and the control constraint sets do not change from one stage to the next. Then, if

$$p_{ij}(u) = P\{x_{k+1} = j \mid x_k = i, u_k = u\}$$

are the transition probabilities, we can alternatively represent the system by the system equation (cf. the discussion of the previous section)

$$x_{k+1} = w_k,$$

where the probability distribution of the disturbance w_k is

$$P\{w_k = j \mid x_k = i, u_k = u\} = p_{ij}(u).$$

Using this system equation and denoting by $g(i, u)$ the expected cost per stage at state i when control u is applied, the DP algorithm can be rewritten as

$$J_k(i) = \min_{u \in U(i)} \left[g(i, u) + E\{J_{k+1}(w_k)\} \right]$$

or equivalently (in view of the distribution of w_k given previously)

$$J_k(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_j p_{ij}(u) J_{k+1}(j) \right].$$

As an illustration, in the machine replacement Example 1.1.3, this algorithm takes the form

$$J_N(i) = 0, \quad i = 1, \dots, n,$$

$$J_k(i) = \min \left[R + g(1) + J_{k+1}(1), g(i) + \sum_{j=i}^n p_{ij} J_{k+1}(j) \right].$$

The two expressions in the above minimization correspond to the two available decisions (replace or not replace the machine).

In the queueing Example 1.1.4, the DP algorithm takes the form

$$J_N(i) = R(i), \quad i = 0, 1, \dots, n,$$

$$J_k(i) = \min \left[r(i) + c_f + \sum_{j=0}^n p_{ij}(u_f) J_{k+1}(j), r(i) + c_s + \sum_{j=0}^n p_{ij}(u_s) J_{k+1}(j) \right].$$

The two expressions in the above minimization correspond to the two possible decisions (fast and slow service).

1.4 STATE AUGMENTATION AND OTHER REFORMULATIONS

We now discuss how to deal with situations where some of the assumptions of the basic problem are violated. Generally, in such cases the problem can be reformulated into the basic problem format. This process is called *state augmentation* because it typically involves the enlargement of the state space. The general guideline in state augmentation is to *include in the enlarged state at time k all the information that is known to the controller at time k and can be used with advantage in selecting u_k* . Unfortunately, state augmentation often comes at a price: the reformulated problem may have very complex state and/or control spaces. We provide some examples.

Time Delays

In many applications the system state x_{k+1} depends not only on the preceding state x_k and control u_k but also on earlier states and controls. In other words, states and controls influence future states with some time delay. Such situations can be handled by state augmentation; the state is expanded to include an appropriate number of earlier states and controls.

For simplicity, assume that there is at most a single period time delay in the state and control; i.e., the system equation has the form

$$x_{k+1} = f_k(x_k, x_{k-1}, u_k, u_{k-1}, w_k), \quad k = 1, 2, \dots, N-1, \quad (1.12)$$

$$x_1 = f_0(x_0, u_0, w_0).$$

Time delays of more than one period can be handled similarly.

If we introduce additional state variables y_k and s_k , and we make the identifications $y_k = x_{k-1}$, $s_k = u_{k-1}$, the system equation (1.12) yields

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \\ s_{k+1} \end{pmatrix} = \begin{pmatrix} f_k(x_k, y_k, u_k, s_k, w_k) \\ x_k \\ u_k \end{pmatrix}. \quad (1.13)$$

By defining $\tilde{x}_k = (x_k, y_k, s_k)$ as the new state, we have

$$\tilde{x}_{k+1} = \tilde{f}_k(\tilde{x}_k, u_k, w_k),$$

where the system function \tilde{f}_k is defined from Eq. (1.13). By using the preceding equation as the system equation and by expressing the cost function in terms of the new state, the problem is reduced to the basic problem without time delays. Naturally, the control u_k should now depend on the new state \tilde{x}_k , or equivalently a policy should consist of functions μ_k of the current state x_k , as well as the preceding state x_{k-1} and the preceding control u_{k-1} .

When the DP algorithm for the reformulated problem is translated in terms of the variables of the original problem, it takes the form

$$\begin{aligned} J_N(x_N) &= g_N(x_N), \\ J_{N-1}(x_{N-1}, x_{N-2}, u_{N-2}) &= \min_{u_{N-1} \in U_{N-1}(x_{N-1})} E_{w_{N-1}} \left\{ g_{N-1}(x_{N-1}, u_{N-1}, w_{N-1}) \right. \\ &\quad \left. + J_N(f_{N-1}(x_{N-1}, x_{N-2}, u_{N-1}, u_{N-2}, w_{N-1})) \right\}, \\ J_k(x_k, x_{k-1}, u_{k-1}) &= \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) \right. \\ &\quad \left. + J_{k+1}(f_k(x_k, x_{k-1}, u_k, u_{k-1}, w_k), x_k, u_k) \right\}, \quad k = 1, \dots, N-2, \\ J_0(x_0) &= \min_{u_0 \in U_0(x_0)} E_{w_0} \left\{ g_0(x_0, u_0, w_0) + J_1(f_0(x_0, u_0, w_0), x_0, u_0) \right\}. \end{aligned}$$

Similar reformulations are possible when time delays appear in the cost; for example, in the case where the cost has the form

$$E \left\{ g_N(x_N, x_{N-1}) + g_0(x_0, u_0, w_0) + \sum_{k=1}^{N-1} g_k(x_k, x_{k-1}, u_k, w_k) \right\}.$$

The extreme case of time delays in the cost arises in the nonadditive form

$$E \{ g_N(x_N, x_{N-1}, \dots, x_0, u_{N-1}, \dots, u_0, w_{N-1}, \dots, w_0) \}.$$

Then, the problem can be reduced to the basic problem format, by taking as augmented state

$$\tilde{x}_k = (x_k, x_{k-1}, \dots, x_0, u_{k-1}, \dots, u_0, w_{k-1}, \dots, w_0)$$

and $E\{g_N(\tilde{x}_N)\}$ as reformulated cost. Policies consist of functions μ_k of the present and past states x_k, \dots, x_0 , the past controls u_{k-1}, \dots, u_0 , and the past disturbances w_{k-1}, \dots, w_0 . Naturally, we must assume that the past disturbances are known to the controller. Otherwise, we are faced with a problem where the state is imprecisely known to the controller. Such problems are known as problems with imperfect state information and will be discussed in Chapter 4.

Correlated Disturbances

Consider the case where the disturbances w_k are correlated over time. A common situation that can be handled efficiently by state augmentation arises when the process w_0, \dots, w_{N-1} can be represented as the output of a linear system driven by independent random variables. As an example, suppose that by using statistical methods, we determine that the evolution of w_k can be modeled by an equation of the form

$$w_k = \lambda w_{k-1} + \xi_k,$$

where λ is a given scalar and $\{\xi_k\}$ is a sequence of independent random vectors with given distribution. Then we can introduce an additional state variable

$$y_k = w_{k-1}$$

and obtain a new system equation

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \begin{pmatrix} f_k(x_k, u_k, \lambda y_k + \xi_k) \\ \lambda y_k + \xi_k \end{pmatrix},$$

where the new state is the pair $\tilde{x}_k = (x_k, y_k)$ and the new disturbance is the vector ξ_k .

More generally, suppose that w_k can be modeled by

$$w_k = C_k y_{k+1},$$

where

$$y_{k+1} = A_k y_k + \xi_k, \quad k = 0, \dots, N-1,$$

A_k, C_k are known matrices of appropriate dimension, and ξ_k are independent random vectors with given distribution (see Fig. 1.4.1). By viewing y_k as an additional state variable, we obtain the new system equation

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \begin{pmatrix} f_k(x_k, u_k, C_k(A_k y_k + \xi_k)) \\ A_k y_k + \xi_k \end{pmatrix}.$$

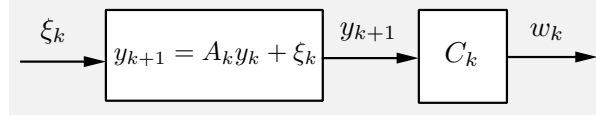


Figure 1.4.1 Representing correlated disturbances as the output of a linear system driven by independent random vectors.

Note that in order to have perfect state information, the controller must be able to observe y_k . Unfortunately, this is true only in the minority of practical cases; for example when C_k is the identity matrix and w_{k-1} is observed before u_k is applied. In the case of perfect state information, the DP algorithm takes the form

$$J_N(x_N, y_N) = g_N(x_N),$$

$$J_k(x_k, y_k) = \min_{u_k \in U_k(x_k)} E_{\xi_k} \left\{ g_k(x_k, u_k, C_k(A_k y_k + \xi_k)) \right. \\ \left. + J_{k+1}(f_k(x_k, u_k, C_k(A_k y_k + \xi_k)), A_k y_k + \xi_k) \right\}.$$

Forecasts

Consider the case where at time k the controller has access to a forecast y_k that results in a reassessment of the probability distribution of w_k and possibly of future disturbances. For example, y_k may be an exact prediction of w_k or an exact prediction that the probability distribution of w_k is a specific one out of a finite collection of distributions. Forecasts of interest in practice are, for example, probabilistic predictions on the state of the weather, the interest rate for money, and the demand for inventory.

Generally, forecasts can be handled by state augmentation although the reformulation into the basic problem format may be quite complex. We will treat here only a simple special case.

Assume that at the beginning of each period k , the controller receives an accurate prediction that the next disturbance w_k will be selected according to a particular probability distribution out of a given collection of distributions $\{Q_1, \dots, Q_m\}$; i.e., if the forecast is i , then w_k is selected according to Q_i . The a priori probability that the forecast will be i is denoted by p_i and is given.

For instance, suppose that in our earlier inventory example the demand w_k is determined according to one of three distributions Q_1 , Q_2 , and Q_3 , corresponding to “small,” “medium,” and “large” demand. Each of the three types of demand occurs with a given probability at each time period, independently of the values of demand at previous time periods. However,

the inventory manager, prior to ordering u_k , gets to know through a forecast the type of demand that will occur. (Note that it is the probability distribution of demand that becomes known through the forecast, not the demand itself.)

The forecasting process can be represented by means of the equation

$$y_{k+1} = \xi_k,$$

where y_{k+1} can take the values $1, \dots, m$, corresponding to the m possible forecasts, and ξ_k is a random variable taking the value i with probability p_i . The interpretation here is that when ξ_k takes the value i , then w_{k+1} will occur according to the distribution Q_i .

By combining the system equation with the forecast equation $y_{k+1} = \xi_k$, we obtain an augmented system given by

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \begin{pmatrix} f_k(x_k, u_k, w_k) \\ \xi_k \end{pmatrix}.$$

The new state is

$$\tilde{x}_k = (x_k, y_k),$$

and because the forecast y_k is known at time k , perfect state information prevails. The new disturbance is

$$\tilde{w}_k = (w_k, \xi_k),$$

and its probability distribution is determined by the distributions Q_i and the probabilities p_i , and depends explicitly on \tilde{x}_k (via y_k) but not on the prior disturbances.

Thus, by suitable reformulation of the cost, the problem can be cast into the basic problem format. Note that the control applied depends on both the current state and the current forecast. The DP algorithm takes the form

$$\begin{aligned} J_N(x_N, y_N) &= g_N(x_N), \\ J_k(x_k, y_k) &= \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) \right. \\ &\quad \left. + \sum_{i=1}^m p_i J_{k+1}(f_k(x_k, u_k, w_k), i) \mid y_k \right\}, \end{aligned} \quad (1.14)$$

where y_k may take the values $1, \dots, m$, and the expectation over w_k is taken with respect to the distribution Q_{y_k} .

It should be clear that the preceding formulation admits several extensions. One example is the case where forecasts can be influenced by the control action and involve several future disturbances. However, the price for these extensions is increased complexity of the corresponding DP algorithm.

Simplification for Uncontrollable State Components

When augmenting the state of a given system we often end up with composite states, consisting of several components. It turns out that if some of these components cannot be affected by the choice of control, the DP algorithm can be simplified considerably, as we will now describe.

Let the state of the system be a composite (x_k, y_k) of two components x_k and y_k . The evolution of the main component, x_k , is affected by the control u_k according to the equation

$$x_{k+1} = f_k(x_k, y_k, u_k, w_k),$$

where the probability distribution $P_k(w_k | x_k, y_k, u_k)$ is given. The evolution of the other component, y_k , is governed by a given conditional distribution $P_k(y_k | x_k)$ and cannot be affected by the control, except indirectly through x_k . One is tempted to view y_k as a disturbance, but there is a difference: y_k is observed by the controller before applying u_k , while w_k occurs after u_k is applied, and indeed w_k may probabilistically depend on u_k .

We will formulate a DP algorithm that is executed over the controllable component of the state, with the dependence on the uncontrollable component being “averaged out.” In particular, let $J_k(x_k, y_k)$ denote the optimal cost-to-go at stage k and state (x_k, y_k) , and define

$$\hat{J}_k(x_k) = E_{y_k} \{ J_k(x_k, y_k) | x_k \}.$$

We will derive a DP algorithm that generates $\hat{J}_k(x_k)$.

Indeed, we have

$$\begin{aligned} \hat{J}_k(x_k) &= E_{y_k} \{ J_k(x_k, y_k) | x_k \} \\ &= E_{y_k} \left\{ \min_{u_k \in U_k(x_k, y_k)} E_{w_k, x_{k+1}, y_{k+1}} \{ g_k(x_k, y_k, u_k, w_k) \right. \\ &\quad \left. + J_{k+1}(x_{k+1}, y_{k+1}) | x_k, y_k, u_k \} | x_k \right\} \\ &= E_{y_k} \left\{ \min_{u_k \in U_k(x_k, y_k)} E_{w_k, x_{k+1}} \{ g_k(x_k, y_k, u_k, w_k) \right. \\ &\quad \left. + E_{y_{k+1}} \{ J_{k+1}(x_{k+1}, y_{k+1}) | x_{k+1} \} | x_k, y_k, u_k \} | x_k \right\}, \end{aligned}$$

and finally

$$\begin{aligned} \hat{J}_k(x_k) &= E_{y_k} \left\{ \min_{u_k \in U_k(x_k, y_k)} E_{w_k} \{ g_k(x_k, y_k, u_k, w_k) \right. \\ &\quad \left. + \hat{J}_{k+1}(f_k(x_k, y_k, u_k, w_k)) \} | x_k \right\}. \end{aligned} \quad (1.15)$$

The advantage of this equivalent DP algorithm is that it is executed over a significantly reduced state space. For example, if x_k takes n possible values and y_k takes m possible values, then DP is executed over n states instead of nm states. Note, however, that the minimization in the right-hand side of the preceding equation yields an optimal control law as a function of the full state (x_k, y_k) .

As an example, consider the augmented state resulting from the incorporation of forecasts, as described earlier in this section. Then, the forecast y_k represents an uncontrolled state component, so that the DP algorithm can be simplified as in Eq. (1.15). In particular, by defining

$$\hat{J}_k(x_k) = \sum_{i=1}^m p_i J_k(x_k, i), \quad k = 0, 1, \dots, N-1,$$

and

$$\hat{J}_N(x_N) = g_N(x_N),$$

we have, using Eq. (1.14),

$$\begin{aligned} \hat{J}_k(x_k) = \sum_{i=1}^m p_i \min_{u_k \in U_k(x_k)} E \Big\{ & g_k(x_k, u_k, w_k) \\ & + \hat{J}_{k+1}(f_k(x_k, u_k, w_k)) \mid y_k = i \Big\}, \end{aligned}$$

which is executed over the space of x_k rather than x_k and y_k . This is a simpler algorithm than the one of Eq. (1.14).

Uncontrollable state components often occur in arrival systems, such as queueing, where action must be taken in response to a random event (such as a customer arrival) that cannot be influenced by the choice of control. Then the state of the arrival system must be augmented to include the random event, but the DP algorithm can be executed over a smaller space, as per Eq. (1.15). Here is another example of similar type.

Example 1.4.1 (Tetris)

Tetris is a popular video game played on a two-dimensional grid. Each square in the grid can be full or empty, making up a “wall of bricks” with “holes” and a “jagged top.” The squares fill up as blocks of different shapes fall from the top of the grid and are added to the top of the wall. As a given block falls, the player can move horizontally and rotate the block in all possible ways, subject to the constraints imposed by the sides of the grid and the top of the wall. The falling blocks are generated independently according to some probability distribution, defined over a finite set of standard shapes. The game starts with an empty grid and ends when a square in the top row becomes full and the top of the wall reaches the top of the grid. When a row of full squares is created, this row is removed, the bricks lying above this

row move one row downward, and the player scores a point. The player's objective is to maximize the score attained (total number of rows removed) within N steps or up to termination of the game, whichever occurs first.

We can model the problem of finding an optimal tetris playing strategy as a stochastic DP problem. The control, denoted by u , is the horizontal positioning and rotation applied to the falling block. The state consists of two components:

- (1) The board position, i.e., a binary description of the full/empty status of each square, denoted by x .
- (2) The shape of the current falling block, denoted by y .

There is also an additional termination state which is cost-free. Once the state reaches the termination state, it stays there with no change in cost.

The shape y is generated according to a probability distribution $p(y)$, independently of the control, so it can be viewed as an uncontrollable state component. The DP algorithm (1.15) is executed over the space of x and has the intuitive form

$$\hat{J}_k(x) = \sum_y p(y) \max_u \left[g(x, y, u) + \hat{J}_{k+1}(f(x, y, u)) \right], \quad \text{for all } x,$$

where $g(x, y, u)$ and $f(x, y, u)$ are the number of points scored (rows removed), and the board position (or termination state) when the state is (x, y) and control u is applied, respectively. Note, however, that despite the simplification in the DP algorithm achieved by eliminating the uncontrollable portion of the state, the number of states x is enormous, and the problem can only be addressed by suboptimal methods, which will be discussed in Chapter 6 and in Vol. II.

1.5 SOME MATHEMATICAL ISSUES

Let us now discuss some technical issues relating to the basic problem formulation and the validity of the DP algorithm. The reader who is not mathematically inclined need not be concerned about these issues and can skip this section without loss of continuity; the mathematical fine points do not contribute significantly to the intuition for solving practical problems and do not matter if the disturbances w_k can take only a finite number of values.

Once an admissible policy $\{\mu_0, \dots, \mu_{N-1}\}$ is adopted, the following sequence of events is envisioned at the typical stage k :

1. The controller observes x_k and applies $u_k = \mu_k(x_k)$.
2. The disturbance w_k is generated according to the given distribution $P_k(\cdot \mid x_k, \mu_k(x_k))$.
3. The cost $g_k(x_k, \mu_k(x_k), w_k)$ is incurred and added to previous costs.

4. The next state x_{k+1} is generated according to the system equation

$$x_{k+1} = f_k(x_k, \mu_k(x_k), w_k).$$

If this is the last stage ($k = N - 1$), the terminal cost $g_N(x_N)$ is added to previous costs and the process terminates. Otherwise, k is incremented, and the same sequence of events is repeated at the next stage.

For each stage, the above process is well-defined and is couched in precise probabilistic terms. Matters are, however, complicated by the need to view the cost as a well-defined random variable with well-defined expected value. The framework of probability theory requires that for each policy we define an underlying probability space, i.e., a set Ω , a collection of events in Ω , and a probability measure on these events. In addition, the cost must be a well-defined random variable on this space in the sense of Appendix C (a measurable function from the probability space into the real line in the terminology of measure-theoretic probability theory). For this to be true, additional (measurability) assumptions on the functions f_k , g_k , and μ_k may be required, and it may be necessary to introduce additional structure on the spaces S_k , C_k , and D_k . Furthermore, these assumptions may restrict the class of admissible policies, since the functions μ_k may be constrained to satisfy additional (measurability) requirements.

Thus, unless these additional assumptions and structure are specified, the basic problem is formulated inadequately from a mathematical point of view. Unfortunately, a rigorous formulation for general state, control, and disturbance spaces is well beyond the mathematical framework of this introductory book and will not be undertaken here. Nonetheless, it turns out that these difficulties are mainly technical and do not substantially affect the basic results to be obtained. For this reason, we find it convenient to proceed with informal derivations and arguments; this is consistent with most of the literature on the subject.

We would like to stress, however, that under at least one frequently satisfied assumption, the mathematical difficulties mentioned above disappear. In particular, let us assume that the disturbance spaces D_k are all countable and the expected values of all terms in the cost are finite for every admissible policy (this is true in particular if the spaces D_k are finite sets). Then, for every admissible policy, the expected values of all the cost terms can be written as (possibly infinite) sums involving the probabilities of the elements of the spaces D_k , and no measurability framework is needed.

Alternatively, one may write the cost as

$$J_\pi(x_0) = E_{x_1, \dots, x_N} \left\{ g_N(x_N) + \sum_{k=0}^{N-1} \tilde{g}_k(x_k, \mu_k(x_k)) \right\}, \quad (1.16)$$

where

$$\tilde{g}_k(x_k, \mu_k(x_k)) = E_{w_k} \left\{ g_k(x_k, \mu_k(x_k), w_k) \mid x_k, \mu_k(x_k) \right\},$$

with the preceding expectation taken with respect to the distribution $P_k(\cdot \mid x_k, \mu_k(x_k))$ defined on the countable set D_k . Then one may take as the basic probability space the Cartesian product of the spaces \tilde{S}_k , $k = 1, \dots, N$, given for all k by

$$\tilde{S}_{k+1} = \{x_{k+1} \in S_{k+1} \mid x_{k+1} = f_k(x_k, \mu_k(x_k), w_k), x_k \in \tilde{S}_k, w_k \in D_k\},$$

where $\tilde{S}_0 = \{x_0\}$. The set \tilde{S}_k is the subset of all states that can be reached at time k when the policy $\{\mu_0, \dots, \mu_{N-1}\}$ is used. Because the disturbance spaces D_k are countable, the sets \tilde{S}_k are also countable (this is true since the union of any countable collection of countable sets is a countable set). The system equation $x_{k+1} = f_k(x_k, \mu_k(x_k), w_k)$, the probability distributions $P_k(\cdot \mid x_k, \mu_k(x_k))$, the initial state x_0 , and the policy $\{\mu_0, \dots, \mu_{N-1}\}$ define a probability distribution on the countable set $\tilde{S}_1 \times \dots \times \tilde{S}_N$, and the expected value in the cost expression (1.16) is defined with respect to this latter distribution.

Let us now give a more detailed proof of the validity of the DP algorithm (Prop. 1.3.1). We assume that the disturbance w_k takes a finite or countable number of values and the expected values of all terms in the expression of the cost function are finite for every admissible policy π . Furthermore, the functions $J_k(x_k)$ generated by the DP algorithm are finite for all states x_k and times k . We do not need to assume that the minimum over u_k in the definition of $J_k(x_k)$ is attained by some $u_k \in U(x_k)$.

For any admissible policy $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$ and each $k = 0, 1, \dots, N-1$, denote $\pi^k = \{\mu_k, \mu_{k+1}, \dots, \mu_{N-1}\}$. For $k = 0, 1, \dots, N-1$, let $J_k^*(x_k)$ be the optimal cost for the $(N-k)$ -stage problem that starts at state x_k and time k , and ends at time N ; i.e.,

$$J_k^*(x_k) = \min_{\pi^k} E \left\{ g_N(x_N) + \sum_{i=k}^{N-1} g_i(x_i, \mu_i(x_i), w_i) \right\}.$$

For $k = N$, we define $J_N^*(x_N) = g_N(x_N)$. We will show by induction that the functions J_k^* are equal to the functions J_k generated by the DP algorithm, so that for $k = 0$, we will obtain the desired result.

For any $\epsilon > 0$, and for all k and x_k , let $\mu_k^\epsilon(x_k)$ attain the minimum in the equation

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} E \{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \}, \quad (1.17)$$

$$k = 0, 1, \dots, N-1,$$

within ϵ ; i.e., for all x_k and k , we have $\mu_k^\epsilon(x_k) \in U_k(x_k)$ and

$$E_{w_k} \{ g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}(f_k(x_k, \mu_k^\epsilon(x_k), w_k)) \} \leq J_k(x_k) + \epsilon. \quad (1.18)$$

Let $J_k^\epsilon(x_k)$ be the expected cost starting at state x_k at time k , and using the policy $\{\mu_k^\epsilon, \mu_{k+1}^\epsilon, \dots, \mu_{N-1}^\epsilon\}$. We will show that for all x_k and k , we have

$$J_k(x_k) \leq J_k^\epsilon(x_k) \leq J_k(x_k) + (N - k)\epsilon, \quad (1.19)$$

$$J_k^*(x_k) \leq J_k^\epsilon(x_k) \leq J_k^*(x_k) + (N - k)\epsilon, \quad (1.20)$$

$$J_k(x_k) = J_k^*(x_k). \quad (1.21)$$

It is seen using Eq. (1.18) that the inequalities (1.19) and (1.20) hold for $k = N - 1$. By taking $\epsilon \rightarrow 0$ in Eqs. (1.19) and (1.20), it is also seen that $J_{N-1} = J_{N-1}^*$. Assume that Eqs. (1.19)-(1.21) hold for index $k + 1$. We will show that they also hold for index k .

Indeed, we have

$$\begin{aligned} J_k^\epsilon(x_k) &= E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}^\epsilon(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} \\ &\leq E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} + (N - k - 1)\epsilon \\ &\leq J_k(x_k) + \epsilon + (N - k - 1)\epsilon \\ &= J_k(x_k) + (N - k)\epsilon, \end{aligned}$$

where the first equation holds by the definition of J_k^ϵ , the first inequality holds by the induction hypothesis, and the second inequality holds Eq. (1.18). We also have

$$\begin{aligned} J_k^\epsilon(x_k) &= E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}^\epsilon(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} \\ &\geq E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} \\ &\geq \min_{u_k \in U(x_k)} E_{w_k} \{g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k))\} \\ &= J_k(x_k), \end{aligned}$$

where the first inequality holds by the induction hypothesis. Combining the preceding two relations, we see that Eq. (1.19) holds for index k .

For every policy $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$, we have

$$\begin{aligned} J_k^\epsilon(x_k) &= E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}^\epsilon(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} \\ &\leq E_{w_k} \{g_k(x_k, \mu_k^\epsilon(x_k), w_k) + J_{k+1}(f_k(x_k, \mu_k^\epsilon(x_k), w_k))\} + (N - k - 1)\epsilon \\ &\leq \min_{u_k \in U(x_k)} E_{w_k} \{g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k))\} + (N - k)\epsilon \\ &\leq \min_{u_k \in U(x_k)} E_{w_k} \{g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k))\} + (N - k)\epsilon \\ &\leq E_{w_k} \{g_k(x_k, \mu_k(x_k), w_k) + J_{\pi, k+1}(f_k(x_k, \mu_k(x_k), w_k))\} + (N - k)\epsilon \\ &= J_{\pi, k}(x_k) + (N - k)\epsilon, \end{aligned}$$

where the first inequality holds by the induction hypothesis, and the second inequality holds by Eq. (1.18). Taking the minimum over π^k in the preceding relation, we obtain for all x_k

$$J_k^\epsilon(x_k) \leq J_k^*(x_k) + (N - k)\epsilon.$$

We also have by the definition of J_k^* , for all x_k ,

$$J_k^*(x_k) \leq J_k^\epsilon(x_k).$$

Combining the preceding two relations, we see that Eq. (1.20) holds for index k . Finally, Eq. (1.21) follows from Eqs. (1.19) and (1.20), by taking $\epsilon \rightarrow 0$, and the induction is complete.

Note that by using $\epsilon = 0$ in the relation

$$J_0^\epsilon(x_k) \leq J_0^*(x_k) + N\epsilon,$$

[cf. Eq. (1.20)], we see that a policy that attains the minimum for all x_k and k in Eq. (1.17) is optimal.

In conclusion, the basic problem has been formulated rigorously, and the DP algorithm has been proved rigorously only when the disturbance spaces D_0, \dots, D_{N-1} are countable sets, and the expected values of all the cost expressions associated with the problem and the DP algorithm are finite. In the absence of these assumptions, the reader should interpret subsequent results and conclusions as essentially correct but mathematically imprecise statements. In fact, when discussing infinite horizon problems (where the need for precision is greater), we will make the countability assumption explicit.

We note, however, that the advanced reader will have little difficulty in establishing most of our subsequent results concerning specific finite horizon applications, even if the countability assumption is not satisfied. This can be done by using the DP algorithm as a verification theorem. In particular, if one can find within a subset of policies $\tilde{\Pi}$ (such as those satisfying certain measurability restrictions) a policy that attains the minimum in the DP algorithm, then this policy can be readily shown to be optimal within $\tilde{\Pi}$. This result is developed in Exercise 1.29, and can be used by the mathematically oriented reader to establish rigorously many of our subsequent results concerning specific applications. For example, in linear-quadratic problems (Section 3.1) one determines from the DP algorithm a policy in closed form, which is linear in the current state. When w_k can take uncountably many values, it is necessary that admissible policies consist of Borel measurable functions μ_k . Since the linear policy obtained from the DP algorithm belongs to this class, the result of Exercise 1.29 guarantees that this policy is optimal.

For a rigorous mathematical treatment of DP that resolves the associated measurability issues and supplements the present text, we refer to the book [BeS78]. Appendix A of Vol. II provides a more accessible survey. The paper [YuB15] describes some recent related developments relating to the policy iteration method (cf. Section 5.3.2).

1.6 DYNAMIC PROGRAMMING AND MINIMAX CONTROL

The problem of optimal control of uncertain systems has traditionally been treated in a stochastic framework, whereby all uncertain quantities are described by probability distributions, and the expected value of the cost is minimized. However, in many practical situations a stochastic description of the uncertainty may not be available, and one may have information with less detailed structure, such as bounds on the magnitude of the uncertain quantities. In other words, one may know a set within which the uncertain quantities are known to lie, but may not know the corresponding probability distribution. Under these circumstances one may use a minimax approach, whereby the worst possible values of the uncertain quantities within the given set are assumed to occur.

The minimax approach for decision making under uncertainty is described in Appendix F and is contrasted with the expected cost approach, which we have been following so far. In its simplest form, the corresponding decision problem is described by a triplet (Π, W, J) , where Π is the set of policies under consideration, W is the set in which the uncertain quantities are known to belong, and $J : \Pi \times W \mapsto [-\infty, +\infty]$ is a given cost function. The objective is to

$$\text{minimize} \quad \max_{w \in W} J(\pi, w)$$

over all $\pi \in \Pi$.

It is possible to formulate a minimax counterpart to the basic problem with perfect state information. This problem is a special case of the abstract minimax problem above, as discussed more fully in Appendix F. Generally, it is unusual for even the simplest special cases of this problem to admit a closed-form solution. However, a computational solution using DP is possible, and our purpose in this section is to describe the corresponding algorithm.

In the framework of the basic problem, consider the case where the disturbances w_0, w_1, \dots, w_{N-1} do not have a probabilistic description but rather are known to belong to corresponding given sets $W_k(x_k, u_k) \subset D_k$, $k = 0, 1, \dots, N-1$, which may depend on the current state x_k and control u_k . Consider the problem of finding a policy $\pi = \{\mu_0, \dots, \mu_{N-1}\}$ with $\mu_k(x_k) \in U_k(x_k)$ for all x_k and k , which minimizes the cost function

$$J_\pi(x_0) = \max_{\substack{w_k \in W_k(x_k, \mu_k(x_k)) \\ k=0,1,\dots,N-1}} \left[g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right].$$

The DP algorithm for this problem takes the following form, which resembles the one corresponding to the stochastic basic problem (maximization is used in place of expectation):

$$J_N(x_N) = g_N(x_N), \tag{1.22}$$

$$J_k(x_k) = \min_{u_k \in U(x_k)} \max_{w_k \in W_k(x_k, u_k)} \left[g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right]. \quad (1.23)$$

This algorithm can be explained by using a principle of optimality type of argument. In particular, we consider the tail subproblem whereby we are at state x_k at time k , and we wish to minimize the “cost-to-go”

$$\max_{\substack{w_i \in W_i(x_i, \mu_i(x_i)) \\ i=k, k+1, \dots, N-1}} \left[g_N(x_N) + \sum_{i=k}^{N-1} g_i(x_i, \mu_i(x_i), w_i) \right],$$

and we argue that if $\pi^* = \{\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*\}$ is an optimal policy for the minimax problem, then the truncated policy $\{\mu_k^*, \mu_{k+1}^*, \dots, \mu_{N-1}^*\}$ is optimal for the tail subproblem. The optimal cost of this subproblem is $J_k(x_k)$, as given by the DP algorithm (1.22)-(1.23). The algorithm expresses the intuitively clear fact that when at state x_k at time k , then regardless of what happened in the past, we should choose u_k that minimizes the worst/maximum value over w_k of the sum of the current stage cost plus the optimal cost of the tail subproblem that starts from the next state.

We will now give a mathematical proof that the DP algorithm (1.22)-(1.23) is valid, and that the optimal cost is equal to $J_0(x_0)$. For this it is necessary to assume that $J_k(x_k) > -\infty$ for all x_k and k . This is analogous to the assumption we made in the preceding section for the validity of the DP algorithm under stochastic disturbances, i.e., that the values $J_k(x_k)$ generated by the DP algorithm are finite for all states x_k and stages k . In the stochastic case the key step of the proof was to bring the minimization over the controls of future stages inside the expectation over the disturbance of the current stage. Similarly, in the minimax case the key step is to bring the minimization over the controls of future stages inside the maximization over the disturbance of the current stage. The following lemma provides the key argument for doing so.

Lemma 1.6.1: Let $f : W \rightarrow X$ be a function, and M be the set of all functions $\mu : X \rightarrow U$, where W , X , and U are some sets. Then for any functions $G_0 : W \rightarrow (-\infty, \infty]$ and $G_1 : X \times U \rightarrow (-\infty, \infty]$ such that

$$\min_{u \in U} G_1(f(w), u) > -\infty, \quad \text{for all } w \in W,$$

we have

$$\min_{\mu \in M} \max_{w \in W} \left[G_0(w) + G_1(f(w), \mu(f(w))) \right] = \max_{w \in W} \left[G_0(w) + \min_{u \in U} G_1(f(w), u) \right].$$

Proof: We have for all $\mu \in M$

$$\max_{w \in W} [G_0(w) + G_1(f(w), \mu(f(w)))] \geq \max_{w \in W} [G_0(w) + \min_{u \in U} G_1(f(w), u)]$$

and by taking the minimum over $\mu \in M$, we obtain

$$\min_{\mu \in M} \max_{w \in W} [G_0(w) + G_1(f(w), \mu(f(w)))] \geq \max_{w \in W} [G_0(w) + \min_{u \in U} G_1(f(w), u)]. \quad (1.24)$$

To show the reverse inequality, for any $\epsilon > 0$, let $\mu_\epsilon \in M$ be such that

$$G_1(f(w), \mu_\epsilon(f(w))) \leq \min_{u \in U} G_1(f(w), u) + \epsilon, \quad \text{for all } w \in W.$$

[Such a μ_ϵ exists because of the assumption $\min_{u \in U} G_1(f(w), u) > -\infty$.]
Then

$$\begin{aligned} \min_{\mu \in M} \max_{w \in W} [G_0(w) + G_1(f(w), \mu(f(w)))] \\ \leq \max_{w \in W} [G_0(w) + G_1(f(w), \mu_\epsilon(f(w)))] \\ \leq \max_{w \in W} [G_0(w) + \min_{u \in U} G_1(f(w), u)] + \epsilon. \end{aligned}$$

Since $\epsilon > 0$ can be taken arbitrarily small, we obtain the reverse to Eq. (1.24), and the desired result follows. **Q.E.D.**

To see how the conclusion of the lemma can fail without the condition

$$\min_{u \in U} G_1(f(w), u) > -\infty$$

for all w , let u be a scalar, let $w = (w_1, w_2)$ be a two-dimensional vector, and let there be no constraints on u and w ($U = \mathbb{R}$, $W = \mathbb{R} \times \mathbb{R}$, where \mathbb{R} is the real line). Let also

$$G_0(w) = w_1, \quad f(w) = w_2, \quad G_1(f(w), u) = f(w) + u.$$

Then, for all $\mu \in M$ we have,

$$\max_{w \in W} [G_0(w) + G_1(f(w), \mu(f(w)))] = \max_{w_1 \in \mathbb{R}, w_2 \in \mathbb{R}} [w_1 + w_2 + \mu(w_2)] = \infty,$$

so that

$$\min_{\mu \in M} \max_{w \in W} [G_0(w) + G_1(f(w), \mu(f(w)))] = \infty.$$

On the other hand,

$$\max_{w \in W} [G_0(w) + \min_{u \in U} G_1(f(w), u)] = \max_{w_1 \in \mathbb{R}, w_2 \in \mathbb{R}} [w_1 + \min_{u \in \mathbb{R}} [w_2 + u]] = -\infty,$$

since $\min_{u \in \mathbb{R}} [w_2 + u] = -\infty$ for all w_2 .

We now turn to proving the DP algorithm (1.22)-(1.23). The proof is similar to the one for the DP algorithm for stochastic problems. The optimal cost $J^*(x_0)$ of the problem is given by

$$\begin{aligned} J^*(x_0) &= \min_{\mu_0} \cdots \min_{\mu_{N-1}} \max_{w_0 \in W[x_0, \mu_0(x_0)]} \cdots \max_{w_{N-1} \in W[x_{N-1}, \mu_{N-1}(x_{N-1})]} \\ &\quad \left[\sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) + g_N(x_N) \right] \\ &= \min_{\mu_0} \cdots \min_{\mu_{N-2}} \left[\min_{\mu_{N-1}} \max_{w_0 \in W[x_0, \mu_0(x_0)]} \cdots \max_{w_{N-2} \in W[x_{N-2}, \mu_{N-2}(x_{N-2})]} \right. \\ &\quad \left[\sum_{k=0}^{N-2} g_k(x_k, \mu_k(x_k), w_k) + \max_{w_{N-1} \in W[x_{N-1}, \mu_{N-1}(x_{N-1})]} \right. \\ &\quad \left. \left. \left[g_{N-1}(x_{N-1}, \mu_{N-1}(x_{N-1}), w_{N-1}) + J_N(x_N) \right] \right] \right]. \end{aligned}$$

We can interchange the minimum over μ_{N-1} and the maximum over w_0, \dots, w_{N-2} by applying Lemma 1.6.1 with the identifications

$$\begin{aligned} w &= (w_0, w_1, \dots, w_{N-2}), \quad u = u_{N-1}, \quad f(w) = x_{N-1}, \\ G_0(w) &= \begin{cases} \sum_{k=0}^{N-2} g_k(x_k, \mu_k(x_k), w_k) & \text{if } w_k \in W_k(x_k, \mu_k(x_k)) \text{ for all } k, \\ \infty & \text{otherwise,} \end{cases} \\ G_1(f(w), u) &= \begin{cases} \hat{G}_1(f(w), u) & \text{if } u \in U_{N-1}(f(w)), \\ \infty & \text{otherwise,} \end{cases} \end{aligned}$$

where

$$\begin{aligned} \hat{G}_1(f(w), u) &= \max_{w_{N-1} \in W_{N-1}(f(w), u)} \left[g_{N-1}(f(w), u, w_{N-1}) \right. \\ &\quad \left. + J_N(f_{N-1}(f(w), u, w_{N-1})) \right], \end{aligned}$$

to obtain

$$\begin{aligned} J^*(x_0) &= \min_{\mu_0} \cdots \min_{\mu_{N-2}} \max_{w_0 \in W[x_0, \mu_0(x_0)]} \cdots \max_{w_{N-2} \in W[x_{N-2}, \mu_{N-2}(x_{N-2})]} \\ &\quad \left[\sum_{k=0}^{N-2} g_k(x_k, \mu_k(x_k), w_k) + J_{N-1}(x_{N-1}) \right]. \end{aligned} \tag{1.25}$$

The required condition $\min_{u \in U} G_1(f(w), u) > -\infty$ for all w (required for application of Lemma 1.6.1) is implied by the assumption $J_{N-1}(x_{N-1}) > -\infty$ for

all x_{N-1} . Now, by working with the expression for $J^*(x_0)$ in Eq. (1.25), and by similarly continuing backwards, with $N - 1$ in place of N , etc., after N steps we obtain

$$J^*(x_0) = J_0(x_0),$$

which is the desired relation. The line of argument just given also shows that an optimal policy for the minimax problem can be constructed by minimizing in the right-hand side of the DP Eq. (1.23), similar to the case of the DP algorithm for the stochastic basic problem.

Unfortunately, as mentioned earlier, there are hardly any interesting examples of an analytical, closed-form solution of the DP algorithm (1.22)-(1.23). A computational solution, requires qualitatively comparable effort to the one of the stochastic DP algorithm. Instead of the expectation operation, one must carry out a maximization operation for each x_k and k .

Minimax control problems will be revisited in Chapter 3 in the context of reachability of target sets and target tubes (Section 3.6.2), and in Chapter 6 in the context of model predictive control (Section 6.4.3).

1.7 NOTES, SOURCES, AND EXERCISES

Dynamic programming is a simple mathematical technique that has been used for many years by engineers, mathematicians, and social scientists in a variety of contexts. It was Bellman, however, who realized in the early fifties that DP could be developed (in conjunction with the then appearing digital computer) into a systematic tool for optimization. In his influential books [Bel57], [BeD62], Bellman demonstrated the broad scope of DP and helped streamline its theory.

Following Bellman's works, the mathematical and algorithmic aspects of infinite horizon problems were extensively investigated, extensions to continuous-time problems were formulated and analyzed, and the mathematical issues discussed in Section 1.5 were addressed. In addition, DP was used in a broad variety of applications, ranging from many branches of engineering to statistics, economics, finance, and some of the social sciences. Samples of these applications will be given in subsequent chapters.

A major methodological advance has been the use of various types of approximations in DP methods for large-scale applications, starting in the late 80s. Considerable success has been obtained in a variety of fields, including prominent achievements with programs that have learned how to play games, such as Backgammon, Go, and others, at impressive and sometimes above human level. We collectively refer to these methods as "approximate DP"; the name "reinforcement learning" is also often used in artificial intelligence, and the names "neuro-dynamic programming" and "adaptive dynamic programming" are often used in automatic control, with essentially the same meaning. We discuss these methods in Chapter 6 and also, more extensively, in Vol. II of this work.

E X E R C I S E S

1.1

Complete the calculations needed to verify that $J_0(1) = 2.7$ and $J_0(2) = 2.818$ in Example 1.3.2.

1.2

Consider the system

$$x_{k+1} = x_k + u_k + w_k, \quad k = 0, 1, 2, 3,$$

with initial state $x_0 = 5$, and the cost function

$$\sum_{k=0}^3 (x_k^2 + u_k^2).$$

Apply the DP algorithm for the following three cases:

- (a) The control constraint set $U_k(x_k)$ is $\{u \mid 0 \leq x_k + u \leq 5, u : \text{integer}\}$ for all x_k and k , and the disturbance w_k is equal to zero for all k .
- (b) The control constraint and the disturbance w_k are as in part (a), but there is in addition a constraint $x_4 = 5$ on the final state. *Hint:* For this problem you need to define a state space for x_4 that consists of just the value $x_4 = 5$, and also to redefine $U_3(x_3)$. Alternatively, you may use a terminal cost $g_4(x_4)$ equal to a very large number for $x_4 \neq 5$.
- (c) The control constraint is as in part (a) and the disturbance w_k takes the values -1 and 1 with equal probability $1/2$ for all x_k and u_k , except if $x_k + u_k$ is equal to 0 or 5 , in which case $w_k = 0$ with probability 1 .

1.3

Suppose we have a machine that is either running or is broken down. If it runs throughout one week, it makes a gross profit of \$100. If it fails during the week, gross profit is zero. If it is running at the start of the week and we perform preventive maintenance, the probability that it will fail during the week is 0.4. If we do not perform such maintenance, the probability of failure is 0.7. However, maintenance will cost \$20. When the machine is broken down at the start of the week, it may either be repaired at a cost of \$40, in which case it will fail during the week with a probability of 0.4, or it may be replaced at a cost of \$150 by a new machine that is guaranteed to run through its first week of operation. Find the optimal repair, replacement, and maintenance policy that maximizes total profit over four weeks, assuming a new machine at the start of the first week.

1.4

A game of the blackjack variety is played by two players as follows: Both players throw a die. The first player, knowing his opponent's result, may stop or may throw the die again and add the result to the result of his previous throw. He then may stop or throw again and add the result of the new throw to the sum of his previous throws. He may repeat this process as many times as he wishes. If his sum exceeds seven (i.e., he busts), he loses the game. If he stops before exceeding seven, the second player takes over and throws the die successively until the sum of his throws is four or higher. If the sum of the second player is over seven, he loses the game. Otherwise the player with the larger sum wins, and in case of a tie the second player wins. The problem is to determine a stopping strategy for the first player that maximizes his probability of winning for each possible initial throw of the second player. Formulate the problem in terms of DP and find an optimal stopping strategy for the case where the second player's initial throw is three. *Hint:* Let $N = 6$ and consider a state space consisting of the following 15 states:

x^1 : busted

x^{1+i} : already stopped at sum i ($1 \leq i \leq 7$),

x^{8+i} : current sum is i but the player has not yet stopped ($1 \leq i \leq 7$).

The optimal strategy is to throw until the sum is four or higher.

1.5 (Computer Assignment)

In the classical game of blackjack the player draws cards knowing only one card of the dealer. The player loses upon reaching a sum of cards exceeding 21. If the player stops before exceeding 21, the dealer draws cards until reaching 17 or higher. The dealer loses upon reaching a sum exceeding 21 or stopping at a lower sum than the player's. If player and dealer end up with an equal sum no one wins. In all other cases the dealer wins. An ace for the player may be counted as a 1 or an 11 as the player chooses. An ace for the dealer is counted as an 11 if this results in a sum from 17 to 21 and as a 1 otherwise. Jacks, queens, and kings count as 10 for both dealer and player. We assume an infinite card deck so the probability of a particular card showing up is independent of earlier cards.

- (a) For every possible initial dealer card, calculate the probability that the dealer will reach a sum of 17, 18, 19, 20, 21, or over 21.
- (b) Calculate the optimal choice of the player (draw or stop) for each of the possible combinations of dealer's card and player's sum of 12 to 20. Assume that the player's cards do not include an ace.
- (c) Repeat part (b) for the case where the player's cards include an ace.

1.6 (Knapsack Problem)

Assume that we have a vessel whose maximum weight capacity is z and whose cargo is to consist of different quantities of N different items. Let v_i denote the value of the i th type of item, w_i the weight of i th type of item, and x_i the number of items of type i that are loaded in the vessel. The problem is to find the most valuable cargo, i.e., to maximize $\sum_{i=1}^N x_i v_i$ subject to the constraints $\sum_{i=1}^N x_i w_i \leq z$ and $x_i = 0, 1, 2, \dots$. Formulate this problem in terms of DP.

1.7 (Traveling Repairman Problem)

A repairman must service n sites, which are located along a line and are sequentially numbered $1, 2, \dots, n$. The repairman starts at a given site s with $1 < s < n$, and is constrained to service only sites that are adjacent to the ones serviced so far, i.e., if he has already serviced sites $i, i+1, \dots, j$, then he may service next only site $i-1$ (assuming $1 < i$) or site $j+1$ (assuming $j < n$). There is a waiting cost c_i for each time period that site i has remained unserved and there is a travel cost t_{ij} for servicing site j immediately after servicing site i . Formulate a DP algorithm for finding a minimum cost service schedule.

1.8 (Ordering Matrix Multiplications) www

Given a sequence of matrix multiplications

$$M_1 M_2 \cdots M_k M_{k+1} \cdots M_N,$$

where each M_k is a matrix of dimension $n_k \times n_{k+1}$, the order in which multiplications are carried out can make a difference. For example, if $n_1 = 1$, $n_2 = 10$, $n_3 = 1$, and $n_4 = 10$, the calculation $((M_1 M_2) M_3)$ requires 20 scalar multiplications, but the calculation $(M_1 (M_2 M_3))$ requires 200 scalar multiplications (multiplying an $m \times n$ matrix with an $n \times k$ matrix requires mnk scalar multiplications).

- (a) Derive a DP algorithm for finding the optimal multiplication order [any order is allowed, including orders that involve multiple partial products each consisting of two or more adjacent matrices, e.g., $((M_1 M_2)(M_3 M_4))$]. Solve the problem for $N = 3$, $n_1 = 2$, $n_2 = 10$, $n_3 = 5$, and $n_4 = 1$.
- (b) Derive a DP algorithm for finding the optimal multiplication order within the class of orders where at each step, we maintain only one partial product that consists only of adjacent matrices, e.g., $((M_1 (M_2 M_3)) M_4)$.

1.9 (Paragraphing Problem)

The paragraphing problem deals with breaking up a sequence of N words of given lengths into lines of length A . Let w_1, \dots, w_N be the words and let L_1, \dots, L_N be their lengths. In a simple version of the problem, words are separated by blanks

whose ideal width is b , but blanks can stretch or shrink if necessary, so that a line $w_i, w_{i+1}, \dots, w_{i+k}$ has length exactly A . The cost associated with the line is $(k+1)|b' - b|$, where $b' = (A - L_i - \dots - L_{i+k})/(k+1)$ is the actual average width of the blanks, except if we have the last line ($N = i+k$), in which case the cost is zero when $b' \geq b$. Formulate a DP algorithm for finding the minimum cost separation. *Hint*: Consider the subproblems of optimally separating w_i, \dots, w_N for $i = 1, \dots, N$.

1.10 (Interval Scheduling)

We have N intervals labeled $1, \dots, N$. The i th interval has start point y_i , end point z_i , and value v_i . We want to select a subset of these intervals that has maximum total value, and such that no pair overlaps. Formulate a DP algorithm to solve this problem. *Hint*: Suppose the intervals are ordered so that $z_1 \leq \dots \leq z_N$. Let the number of periods be N and the states be z_1, \dots, z_N . Let also the optimal value at state z_k be the maximal value over nonoverlapping intervals whose start time is greater than z_i .

1.11

Consider a smaller version of a popular puzzle game. Three square tiles numbered 1, 2, and 3 are placed in a 2×2 grid with one space left empty. The two tiles adjacent to the empty space can be moved into that space, thereby creating new configurations. Use a DP argument to answer the question whether it is possible to generate a given configuration starting from any other configuration.

1.12

From a pile of eleven matchsticks, two players take turns removing one or four sticks. The player who removes the last stick wins. Use a DP argument to show that there is a winning strategy for the player who plays first.

1.13 (Counterfeit Coin Problem)

We are given six coins, one of which is counterfeit and is known to have different weight than the rest. Construct a strategy to find the counterfeit coin using a two-pan scale in a minimum average number of tries. *Hint*: There are two initial decisions that make sense: (1) test two of the coins against two others, and (2) test one of the coins against one other.

1.14 (Multiplicative Cost)

In the framework of the basic problem, consider the case where the cost has the multiplicative form

$$\prod_{k=0,1,\dots,N-1}^{E_{w_k}} \{g_N(x_N) \cdot g_{N-1}(x_{N-1}, u_{N-1}, w_{N-1}) \cdots g_0(x_0, u_0, w_0)\}.$$

Develop a DP-like algorithm for this problem assuming that $g_k(x_k, u_k, w_k) \geq 0$ for all x_k, u_k, w_k , and k .

1.15

Consider a device consisting of N stages connected in series, where each stage consists of a particular component. The components are subject to failure, and to increase the reliability of the device duplicate components are provided. For $j = 1, 2, \dots, N$, let $(1 + m_j)$ be the number of components for the j th stage, let $p_j(m_j)$ be the probability of successful operation of the j th stage when $(1 + m_j)$ components are used, and let c_j denote the cost of a single component at the j th stage. Formulate in terms of DP the problem of finding the number of components at each stage that maximize the reliability of the device expressed by

$$p_1(m_1) \cdot p_2(m_2) \cdots p_N(m_N),$$

subject to the cost constraint $\sum_{j=1}^N c_j m_j \leq A$, where $A > 0$ is given.

1.16 www

An innkeeper charges a different rate for a room as the day progresses, depending on whether he has many or few vacancies. His objective is to maximize his expected total income during the day. Let x be the number of empty rooms at the start of the day, and let y be the number of customers that will ask for a room in the course of the day. We assume (somewhat unrealistically) that the innkeeper knows y with certainty, and upon arrival of a customer, quotes one of m prices r_i , $i = 1, \dots, m$, where $0 < r_1 \leq r_2 \leq \dots \leq r_m$. A quote of a rate r_i is accepted with probability p_i and is rejected with probability $1 - p_i$, in which case the customer departs, never to return during that day.

- (a) Formulate this as a problem with y stages and show that the maximal expected income, as a function of x and y , satisfies the recursion

$$J(x, y) = \max_{i=1, \dots, m} \left[p_i(r_i + J(x-1, y-1)) + (1 - p_i)J(x, y-1) \right],$$

for all $x \geq 1$ and $y \geq 1$, with initial conditions

$$J(x, 0) = J(0, y) = 0, \quad \text{for all } x \text{ and } y.$$

Assuming that the product $p_i r_i$ is monotonically nondecreasing with i , and that p_i is monotonically nonincreasing with i , show that the innkeeper should always charge the highest rate r_m .

- (b) Consider a variant of the problem where each arriving customer, with probability p_i , offers a price r_i for a room, which the innkeeper may accept or reject. In the latter case the customer departs, never to return during that day. Show that an appropriate DP algorithm is

$$J(x, y) = \sum_{i=1}^m p_i \max[r_i + J(x-1, y-1), J(x, y-1)],$$

with initial conditions

$$J(x, 0) = J(0, y) = 0, \quad \text{for all } x \text{ and } y.$$

Show also that for given x and y it is optimal to accept a customer's offer if it is larger than some threshold $\bar{\tau}(x, y)$. *Hint:* This part is related to DP for uncontrollable state components (cf. Section 1.4).

1.17 (Investing in a Stock) www

An investor observes at the beginning of each period k the price x_k of a stock and decides whether to buy 1 unit, sell 1 unit, or do nothing. There is a transaction cost c for buying or selling. The stock price can take one of n different values v^1, \dots, v^n and the transition probabilities

$$p_{ij}^k = P\{x_{k+1} = v^j \mid x_k = v^i\}$$

are known. The investor wants to maximize the total worth of his stock at a fixed final period N minus his investment costs from period 0 to period $N - 1$ (revenue from a sale is viewed as negative cost). We assume that the function

$$P_k(x) = E\{x_N \mid x_k = x\} - x$$

is monotonically nonincreasing as a function of x ; i.e., the expected profit from a purchase is a nonincreasing function of the purchase price.

- (a) Assume that the investor starts with N or more units of stock and an unlimited amount of cash, so that a purchase or sale decision is possible at each period regardless of the past decisions and the current price. For every period k , let \underline{x}_k be the largest value of $x \in \{v^1, \dots, v^n\}$ such that $P_k(x) > c$, and let \bar{x}_k be the smallest value of $x \in \{v^1, \dots, v^n\}$ such that $P_k(x) < -c$. Show that it is optimal to buy if $x_k \leq \underline{x}_k$, sell if $\bar{x}_k \leq x_k$, and do nothing otherwise. *Hint:* Formulate the problem as one of maximizing

$$E \left\{ \sum_{k=0}^{N-1} (u_k P_k(x_k) - c |u_k|) \right\},$$

where $u_k \in \{-1, 0, 1\}$.

- (b) Formulate an efficient DP algorithm for the case where the investor starts with less than N units of stock and an unlimited amount of cash. Show that it is still optimal to buy if $x_k \leq \underline{x}_k$ and it is still not optimal to sell if $x_k < \bar{x}_k$. Could it be optimal to buy at any prices x_k greater than \underline{x}_k ?
- (c) Consider the situation where the investor initially has N or more units of stock and there is a constraint that for any time k the number of purchase decisions up to k should not exceed the number of sale decisions up to k by more than a given fixed number m (this models approximately the situation where the investor has a limited initial amount of cash). Formulate an

efficient DP algorithm for this case. Show that it is still optimal to sell if $\bar{x}_k \leq x_k$ and it is still not optimal to buy if $\underline{x}_k < x_k$.

- (d) Consider the situation where there are restrictions on both the initial amount of stock as in part (b), and the number of purchase decisions as in part (c). Derive a DP algorithm for this problem.
- (e) How would the analysis of (a)-(d) be affected if cash is invested at a given fixed interest rate?

1.18 (Regular Polygon Theorem) www

According to a famous theorem (attributed to the ancient Greek geometer Zenodorus), of all N -side polygons inscribed in a given circle, those that are regular (all sides are equal) have maximal area.

- (a) Prove the theorem by applying DP to a suitable problem involving sequential placement of N points in the circle.
- (b) Use DP to solve the problem of placing a given number of points on a subarc of the circle, so as to maximize the area of the polygon whose vertices are these points, the endpoints of the subarc, and the center of the circle.

1.19 (Inscribed Polygon of Maximal Perimeter)

Consider the problem of inscribing an N -side polygon in a given circle, so that the polygon has maximal perimeter.

- (a) Formulate the problem as a DP problem involving sequential placement of N points in the circle.
- (b) Use DP to show that the optimal polygon is regular (all sides are equal).

1.20 www

A decision maker must continually choose between two activities over a time interval $[0, T]$. Choosing activity i at time t , where $i = 1, 2$, earns reward at a rate $g_i(t)$, and every switch between the two activities costs $c > 0$. Thus, for example, the reward for starting with activity 1, switching to 2 at time t_1 , and switching back to 1 at time $t_2 > t_1$ earns total reward

$$\int_0^{t_1} g_1(t) dt + \int_{t_1}^{t_2} g_2(t) dt + \int_{t_2}^T g_1(t) dt - 2c.$$

We want to find a set of switching times that maximize the total reward. Assume that the function $g_1(t) - g_2(t)$ changes sign a finite number of times in the interval $[0, T]$. Formulate the problem as a finite horizon problem and write the corresponding DP algorithm. See Shreve [Shr81] for a fuller development of this problem.

1.21

A farmer annually producing x_k units of a certain crop stores $(1 - u_k)x_k$ units of his production, where $0 \leq u_k \leq 1$, and invests the remaining $u_k x_k$ units, thus increasing the next year's production to a level x_{k+1} given by

$$x_{k+1} = x_k + w_k u_k x_k, \quad k = 0, 1, \dots, N-1.$$

The scalars w_k are independent random variables with identical probability distributions that do not depend either on x_k or u_k . Furthermore, $E\{w_k\} = \bar{w} > 0$. The problem is to find the optimal investment policy that maximizes the total expected product stored over N years

$$E_{w_k, k=0,1,\dots,N-1} \left\{ x_N + \sum_{k=0}^{N-1} (1 - u_k) x_k \right\}.$$

Show the optimality of the following policy that consists of constant functions:

- (a) If $\bar{w} > 1$, $\mu_0^*(x_0) = \dots = \mu_{N-1}^*(x_{N-1}) = 1$.
- (b) If $0 < \bar{w} < 1/N$, $\mu_0^*(x_0) = \dots = \mu_{N-1}^*(x_{N-1}) = 0$.
- (c) If $1/N \leq \bar{w} \leq 1$,

$$\mu_0^*(x_0) = \dots = \mu_{N-\bar{k}-1}^*(x_{N-\bar{k}-1}) = 1,$$

$$\mu_{N-\bar{k}}^*(x_{N-\bar{k}}) = \dots = \mu_{N-1}^*(x_{N-1}) = 0,$$

where \bar{k} is such that $1/(\bar{k} + 1) < \bar{w} \leq 1/\bar{k}$.

1.22

Let x_k denote the number of educators in a certain country at time k and let y_k denote the number of research scientists at time k . New scientists (potential educators or research scientists) are produced during the k th period by educators at a rate γ_k per educator, while educators and research scientists leave the field due to death, retirement, and transfer at a rate δ_k . The scalars γ_k , $k = 0, 1, \dots, N-1$, are independent identically distributed random variables taking values within a closed and bounded interval of positive numbers. Similarly δ_k , $k = 0, 1, \dots, N-1$, are independent identically distributed and take values in an interval $[\delta, \delta']$ with $0 < \delta \leq \delta' < 1$. By means of incentives, a science policy maker can determine the proportion u_k of new scientists produced at time k who become educators. Thus, the number of research scientists and educators evolves according to the equations

$$x_{k+1} = (1 - \delta_k)x_k + u_k \gamma_k x_k,$$

$$y_{k+1} = (1 - \delta_k)y_k + (1 - u_k)\gamma_k x_k.$$

The initial numbers x_0, y_0 are known, and it is required to find a policy

$$\{\mu_0^*(x_0, y_0), \dots, \mu_{N-1}^*(x_{N-1}, y_{N-1})\}$$

with

$$0 < \alpha \leq \mu_k^*(x_k, y_k) \leq \beta < 1, \quad \text{for all } x_k, y_k, \text{ and } k,$$

which maximizes $E_{\gamma_k, \delta_k} \{y_N\}$ (i.e., the expected final number of research scientists after N periods). The scalars α and β are given.

- (a) Show that the cost-to-go functions $J_k(x_k, y_k)$ are linear; i.e., for some scalars ξ_k, ζ_k ,

$$J_k(x_k, y_k) = \xi_k x_k + \zeta_k y_k.$$

- (b) Derive an optimal policy $\{\mu_0^*, \dots, \mu_{N-1}^*\}$ under the assumption

$$E\{\gamma_k\} > E\{\delta_k\}$$

and show that this optimal policy can consist of constant functions.

- (c) Assume that the proportion of new scientists who become educators at time k is $u_k + \epsilon_k$ (rather than u_k), where ϵ_k are identically distributed independent random variables that are also independent of γ_k, δ_k and take values in the interval $[-\alpha, 1-\beta]$. Derive the form of the cost-to-go functions and the optimal policy.

1.23 (Discounted Cost per Stage)

In the framework of the basic problem, consider the case where the cost is of the form

$$E_{w_k} \left\{ \alpha^N g_N(x_N) + \sum_{k=0}^{N-1} \alpha^k g_k(x_k, u_k, w_k) \right\},$$

where α is a discount factor with $0 < \alpha < 1$. Show that an alternate form of the DP algorithm is given by

$$V_N(x_N) = g_N(x_N),$$

$$V_k(x_k) = \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) + \alpha V_{k+1}(f_k(x_k, u_k, w_k)) \right\}.$$

1.24 (Exponential Cost Function)

In the framework of the basic problem, consider the case where the cost is of the form

$$E_{w_k} \left\{ \exp \left(g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right) \right\}.$$

- (a) Show that the optimal cost and an optimal policy can be obtained from the DP-like algorithm

$$J_N(x_N) = \exp(g_N(x_N)),$$

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ J_{k+1}(f_k(x_k, u_k, w_k)) \exp(g_k(x_k, u_k, w_k)) \right\}.$$

- (b) Define the functions $V_k(x_k) = \ln J_k(x_k)$. Assume also that g_k is a function of x_k and u_k only (and not of w_k). Show that the above algorithm can be rewritten as

$$V_N(x_N) = g_N(x_N),$$

$$V_k(x_k) = \min_{u_k \in U_k(x_k)} \left\{ g_k(x_k, u_k) + \ln E_{w_k} \left\{ \exp(V_{k+1}(f_k(x_k, u_k, w_k))) \right\} \right\}.$$

Note: The exponential is an example of a *risk-sensitive cost function* that can be used to encode a preference for policies with a small variance of the cost $g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k)$. The associated problems have a lot of interesting properties, which are discussed in several sources, e.g., [DeR79], [Whi90], [FeM94], [JBE94], [BaB95], [Bas00], [Pat01], [Ber16b].

1.25 (Terminating Process)

In the framework of the basic problem, consider the case where the system evolution terminates at time i when a given value \bar{w}_i of the disturbance at time i occurs, or when a termination decision u_i is made by the controller. If termination occurs at time i , the resulting cost is

$$T + \sum_{k=0}^i g_k(x_k, u_k, w_k),$$

where T is a termination cost. If the process has not terminated up to the final time N , the resulting cost is $g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k)$. Reformulate the problem into the framework of the basic problem. *Hint:* Augment the state space with a special termination state.

1.26

Alexei plays a game that starts with a deck with b “black” cards and r “red” cards. At each time period he draws a random card and decides between the following two options:

- (1) Without looking at the card, “predict” that it is black, in which case he wins the game if the prediction is correct and loses if the prediction is incorrect.
- (2) “Discard” the card, after looking at its color, and continue the game with one card less.

If the deck has only black cards he wins the game, while if the deck has only red cards he loses the game. Alexei wants to find a policy that maximizes his probability of a win.

- (a) Formulate Alexei's problem into the format of the finite-horizon basic problem with perfect state information. Identify states, controls, and disturbances, and write the DP algorithm.
- (b) Use induction to show that the optimal probability of a win starting with b black cards and r red cards is $\frac{b}{b+r}$.
- (c) Characterize the optimal policies.
- (d) Suppose that Alexei is given the additional option to randomize his decision at each time period. In particular, he may choose a probability $p \in [0, 1]$, flip a coin that has probability of head equal to p , and decide upon option 1 or 2 above depending on the outcome of the flip. What would then be the optimal policies?

1.27 (Semilinear Systems) www

Consider a problem involving the system

$$x_{k+1} = A_k x_k + f_k(u_k) + w_k,$$

where $x_k \in \mathbb{R}^n$, f_k are given functions, and A_k and w_k are random $n \times n$ matrices and n -vectors, respectively, with given probability distributions that do not depend on x_k , u_k or prior values of A_k and w_k . Assume that the cost function is linear in the states and has the form

$$E_{\substack{A_k, w_k \\ k=0,1,\dots,N-1}} \left\{ c'_N x_N + \sum_{k=0}^{N-1} \left(c'_k x_k + g_k(\mu_k(x_k)) \right) \right\},$$

where c_k are given vectors and g_k are given functions. Show that if the optimal cost for this problem is finite and the control constraint sets $U_k(x_k)$ are independent of x_k , then the cost-to-go functions of the DP algorithm are affine (linear plus constant). Assuming that there is at least one optimal policy, show that there exists an optimal policy that is open-loop, i.e., $\mu_k^*(x_k) = \text{constant}$ for all $x_k \in \mathbb{R}^n$.

1.28 (Monotonicity Property of DP) www

An evident, yet very important property of the DP algorithm is that if the terminal cost g_N is changed to a uniformly larger cost \bar{g}_N [i.e., $g_N(x_N) \leq \bar{g}_N(x_N)$ for all x_N], then the last stage cost-to-go $J_{N-1}(x_{N-1})$ will be uniformly increased. More generally, given two functions J_{k+1} and \bar{J}_{k+1} with $J_{k+1}(x_{k+1}) \leq \bar{J}_{k+1}(x_{k+1})$ for all x_{k+1} , we have, for all x_k and $u_k \in U_k(x_k)$,

$$\begin{aligned} E_{w_k} \left\{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right\} \\ \leq E_{w_k} \left\{ g_k(x_k, u_k, w_k) + \bar{J}_{k+1}(f_k(x_k, u_k, w_k)) \right\}. \end{aligned}$$

Suppose now that in the basic problem the system and cost are time invariant; i.e., $S_k \equiv S$, $C_k \equiv C$, $D_k \equiv D$, $f_k \equiv f$, $U_k \equiv U$, $P_k \equiv P$, and $g_k \equiv g$ for some S , C , D , f , U , P , and g . Use induction to show that if in the DP algorithm we have $J_{N-1}(x) \leq J_N(x)$ for all $x \in S$, then

$$J_k(x) \leq J_{k+1}(x), \quad \text{for all } x \in S \text{ and } k.$$

Similarly, if we have $J_{N-1}(x) \geq J_N(x)$ for all $x \in S$, then

$$J_k(x) \geq J_{k+1}(x), \quad \text{for all } x \in S \text{ and } k.$$

1.29 (DP Algorithm for Minimization over a Subset of Policies)

This exercise is primarily of theoretical interest (see the discussion at the end of Section 1.5), but also relates to situations where we can guess that an optimal policy may be found within a special class of policies. Consider a variation of the basic problem whereby we want to find

$$\min_{\pi \in \tilde{\Pi}} J_\pi(x_0),$$

where $\tilde{\Pi}$ is some given *subset* of the set of sequences $\{\mu_0, \mu_1, \dots, \mu_{N-1}\}$ of functions $\mu_k : S_k \rightarrow C_k$ with $\mu_k(x_k) \in U_k(x_k)$ for all $x_k \in S_k$. Assume that for every $\pi = \{\mu_0, \dots, \mu_{N-1}\} \in \tilde{\Pi}$, the sequence of cost-to-go functions $\tilde{J}_{\pi,k}$, $k = 0, \dots, N$, generated by

$$\tilde{J}_{\pi,N}(x_N) = g_N(x_N),$$

$$\tilde{J}_{\pi,k}(x_k) = E_{w_k} \left\{ g_k(x_k, \mu_k(x_k), w_k) + \tilde{J}_{\pi,k+1}(f_k(x_k, \mu_k(x_k), w_k)) \right\},$$

is well-defined in the sense that the functions $\tilde{J}_{\pi,k}$ are real-valued, and that the expected value in the preceding equation is well-defined and finite. Suppose also that

$$\tilde{\pi} = \{\tilde{\mu}_0, \tilde{\mu}_1, \dots, \tilde{\mu}_{N-1}\}$$

is a policy that belongs to $\tilde{\Pi}$ and attains the minimum in the DP algorithm, in the sense that for all x_k and $k = 0, \dots, N-1$, we have

$$\begin{aligned} & E_{w_k} \left\{ g_k(x_k, \tilde{\mu}_k(x_k), w_k) + \tilde{J}_{\tilde{\pi},k+1}(f_k(x_k, \tilde{\mu}_k(x_k), w_k)) \right\} \\ &= \min_{u_k \in U_k(x_k)} E_{w_k} \left\{ g_k(x_k, u_k, w_k) + \tilde{J}_{\tilde{\pi},k+1}(f_k(x_k, u_k, w_k)) \right\}. \end{aligned}$$

Show that $\tilde{\pi}$ is optimal within $\tilde{\Pi}$ in the sense that $J_{\tilde{\pi},0}(x_0) \leq J_{\pi,0}(x_0)$ for all $\pi \in \tilde{\Pi}$ and states x_0 . *Hint:* Use backwards induction to show that $J_{\tilde{\pi},k}(x_k) \leq J_{\pi,k}(x_k)$ for all $\pi \in \tilde{\Pi}$, k , and states x_k .

1.30 (Post-Decision States)

Consider the basic problem and assume that the system equation has a special structure whereby from state x_k after applying u_k we move to an intermediate “post-decision state” $y_k = p_k(x_k, u_k)$ at cost $g_k(x_k, u_k)$. Then from y_k we move at no cost to the new state x_{k+1} according to

$$x_{k+1} = h_k(y_k, w_k),$$

where the distribution of the disturbance w_k depends only on y_k , and not on prior disturbances, states, and controls. The purpose of this exercise is to show that it is possible to exploit the structure of the problem to execute the DP algorithm more efficiently. Denote by $J_k(x_k)$ the optimal cost-to-go starting at time k from state x_k , and by $V_k(y_k)$ the optimal cost-to-go starting at time k from post-decision state y_k .

(a) Show that a DP algorithm that generates only J_k is given by

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} [g(x_k, u_k) + E_{w_k} \{ J_{k+1}(h_k(p_k(x_k, u_k), w_k)) \}].$$

(b) Show that a DP algorithm that generates both J_k and V_k is given by

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} [g(x_k, u_k) + V_k(p_k(x_k, u_k))],$$

$$V_k(y_k) = E_{w_k} \{ J_{k+1}(h_k(y_k, w_k)) \}.$$

(c) Show that a DP algorithm that generates only V_k for all k is given by

$$V_k(y_k) = E_{w_k} \left\{ \min_{u_{k+1} \in U_{k+1}(h_k(y_k, w_k))} [g_{k+1}(h_k(y_k, w_k), u_{k+1}) + V_{k+1}(p_{k+1}(h_k(y_k, w_k), u_{k+1}))] \right\}.$$

5

Introduction to Infinite Horizon Problems

Contents

5.1. An Overview	p. 232
5.2. Stochastic Shortest Path Problems	p. 236
5.3. Computational Methods	p. 245
5.3.1. Value Iteration	p. 245
5.3.2. Policy Iteration	p. 246
5.3.3. Linear Programming	p. 248
5.4. Discounted Problems	p. 249
5.5. Average Cost per Stage Problems	p. 253
5.6. Semi-Markov Problems	p. 267
5.7. Notes, Sources, and Exercises	p. 277

In this chapter, we provide an introduction to infinite horizon problems. These problems differ from those considered so far in two respects:

- (a) The number of stages is infinite.
- (b) The system is stationary, i.e., the system equation, the cost per stage, and the random disturbance statistics do not change from one stage to the next.

The assumption of an infinite number of stages is never satisfied in practice, but is a reasonable approximation for problems involving a finite but very large number of stages. The assumption of stationarity is often satisfied in practice, and in other cases it approximates well a situation where the system parameters vary relatively slowly with time.

Infinite horizon problems are interesting in that their analysis is elegant and insightful, and the implementation of optimal policies is often simple. For example, optimal policies are typically stationary, i.e., the optimal rule for choosing controls does not change from one stage to the next.

On the other hand, infinite horizon problems generally require more sophisticated analysis than their finite horizon counterparts, because of the need to analyze limiting behavior as the horizon tends to infinity. This analysis is often nontrivial and at times reveals surprising possibilities. Our treatment will be limited to finite-state problems. A far more detailed development, together with applications from a variety of fields, and an extensive discussion of the associated suboptimal control and DP approximation issues can be found in Vol. II of this work.

5.1 AN OVERVIEW

There are four principal classes of infinite horizon problems. In the first three classes, we try to minimize the *total cost over an infinite number of stages*, given by

$$J_{\pi}(x_0) = \lim_{N \rightarrow \infty} E_{w_k} \left\{ \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\}.$$

(We assume that the limit above exists for the moment, and address the issue later.) Here, $J_{\pi}(x_0)$ denotes the cost associated with an initial state x_0 and a policy $\pi = \{\mu_0, \mu_1, \dots\}$, and α is a positive scalar with $0 < \alpha \leq 1$, called the *discount factor*. The meaning of $\alpha < 1$ is that future costs matter to us less than the same costs incurred at the present time. As an example, think of k th period dollars depreciated to initial period dollars by a factor of $(1+r)^{-k}$, where r is a rate of interest; here $\alpha = 1/(1+r)$. An important concern in total cost problems is that the limit in the definition

of $J_\pi(x_0)$ be finite. In the first two of the following classes of problems, this is guaranteed through various assumptions on the problem structure and the discount factor. In the third class, the analysis is adjusted to deal with infinite cost for some of the policies. In the fourth class, this cost need not be finite for any policy, and for this reason, the cost is appropriately redefined.

- (a) *Stochastic shortest path problems.* Here, $\alpha = 1$ but there is a special cost-free termination state; once the system reaches that state it remains there at no further cost. We will assume a problem structure such that termination is inevitable (this assumption will be relaxed somewhat in Chapters 3 and 4 of Vol. II). Thus the horizon is in effect finite, but its length is random and may be affected by the policy being used. These problems will be considered in the next section and their analysis will provide the foundation for the analysis of the other types of problems considered in this chapter.
- (b) *Discounted problems with bounded cost per stage.* Here, $\alpha < 1$ and the absolute cost per stage $|g(x, u, w)|$ is bounded from above by some constant M ; this makes the limit in the definition of the cost $J_\pi(x_0)$ well defined because it is the infinite sum of a sequence of numbers that are bounded in absolute value by the decreasing geometric progression $\{\alpha^k M\}$. We will consider these problems in Section 5.4.
- (c) *Discounted and undiscounted problems with unbounded cost per stage.* Here the discount factor α may or may not be less than 1, and the cost per stage may be unbounded. These problems require a complicated analysis because the possibility of infinite cost for some of the policies is explicitly dealt with. We will not consider these problems here; see Chapter 4 of Vol. II.
- (d) *Average cost per stage problems.* Minimization of the total cost $J_\pi(x_0)$ makes sense only if $J_\pi(x_0)$ is finite for at least some admissible policies π and some initial states x_0 . Frequently, however, it turns out that $J_\pi(x_0) = \infty$ for every policy π and initial state x_0 (think of the case where $\alpha = 1$, and the cost for every state and control is positive). It turns out that in many such problems the *average cost per stage*, defined by

$$\lim_{N \rightarrow \infty} \frac{1}{N} E_{w_k} \left\{ \sum_{k=0}^{N-1} g(x_k, \mu_k(x_k), w_k) \right\},$$

is well defined and finite. We will consider some of these problems in Section 5.5.

A Preview of Total Cost Infinite Horizon Results

There are several analytical and computational issues regarding total cost

infinite horizon problems. Many of these revolve around the relation between the optimal cost-to-go function J^* of the infinite horizon problem and the optimal cost-to-go functions of the corresponding N -stage problems. In particular, consider the case $\alpha = 1$ and let $J_N(x)$ denote the optimal cost of the problem involving N stages, initial state x , cost per stage $g(x, u, w)$, and zero terminal cost. The optimal N -stage cost is generated after N iterations of the DP algorithm

$$J_{k+1}(x) = \min_{u \in U(x)} E_w \left\{ g(x, u, w) + J_k(f(x, u, w)) \right\}, \quad k = 0, 1, \dots, \quad (5.1)$$

starting from the initial condition $J_0(x) = 0$ for all x .[†] Since the infinite horizon cost of a given policy is, by definition, the limit of the corresponding N -stage costs as $N \rightarrow \infty$, it is natural to speculate that:

- (1) The optimal infinite horizon cost is the limit of the corresponding N -stage optimal costs as $N \rightarrow \infty$; i.e.,

$$J^*(x) = \lim_{N \rightarrow \infty} J_N(x) \quad (5.2)$$

for all states x . This relation is extremely valuable computationally and analytically, and, fortunately, it typically holds. In particular, it holds for the models of the next two sections [categories (a) and (b) above]. However, there are some unusual exceptions for problems in category (c) above, and this illustrates that infinite horizon problems should be approached with some care. This issue is discussed in more detail in Vol. II.

- (2) The following limiting form of the DP algorithm should hold for all states x ,

$$J^*(x) = \min_{u \in U(x)} E_w \left\{ g(x, u, w) + J^*(f(x, u, w)) \right\},$$

as suggested by Eqs. (5.1) and (5.2). This is not really an algorithm, but rather a system of equations (one equation per state), which has as solution the costs-to-go of all the states. It can also be viewed as a *functional equation* for the cost-to-go function J^* , and it is called

[†] Note here that we have reversed the time indexing to suit our purposes. Also in the case where there is a discount factor $\alpha < 1$, the appropriate form of Eq. (5.1) is

$$J_{k+1}(x) = \min_{u \in U(x)} E_w \left\{ g(x, u, w) + \alpha J_k(f(x, u, w)) \right\}, \quad k = 0, 1, \dots,$$

as will be seen in Section 5.4.

Bellman's equation. Fortunately again, an appropriate form of this equation holds for every type of infinite horizon problem of interest [with some rare exceptions for some problems in categories (a) and (c); see Chapter 4 of Vol. II].

- (3) If $\mu(x)$ attains the minimum in the right-hand side of Bellman's equation for each x , then the policy $\{\mu, \mu, \dots\}$ should be optimal. This is true for most infinite horizon problems of interest and in particular, for all the models discussed in this chapter.

Most of the analysis of infinite horizon problems revolves around the above three issues and also around the issue of efficient computation of J^* and an optimal policy. In the next three sections we will provide a discussion of these issues for some of the simpler infinite horizon problems, all of which involve a finite state space. In Sections 5.5 and 5.6, we discuss extensions to average cost and other problems.

From a mathematical point of view, the fundamental property of a total cost problem (which determines to a large extent the nature of the results that one can prove about it) is whether the DP mapping that transforms one cost-to-go function into another is a *contraction mapping* in a mathematical sense. We postpone this type of mathematical view of infinite horizon problems for Vol. II, but to provide a high-level summary, we note that total cost problems may also be divided into:

- (1) *Contractive* problems, where for all policies the DP mapping is a contraction. Discounted problems with bounded cost per stage [category (b)] form the main class of problems that are contractive, and mathematically speaking, this is the fundamental reason for their good behavior.
- (2) *Noncontractive* problems, where the contraction property just noted does not hold.
- (3) *Semicontractive* problems, where the contraction property holds for some policies, which we refer to as “regular” but not for others. In the most well-behaved type of such problems, there are enough assumptions that guarantee that the search for an optimal policy may be restricted to within the “regular” class. Many stochastic shortest path problems belong to this category.

This mathematical line of analysis has the advantage that it provides a unifying view of infinite horizon DP, which is uncluttered by specialized assumptions, and applies to broader classes of problems. We refer to Vol. II, Chapters 3 and 4, and to the monograph [Ber13a] and subsequent paper [Ber15b] for further discussion.

Total Cost Problem Formulation

Throughout this chapter we assume a controlled finite-state discrete-time

dynamic system, and we will use the corresponding transition probability notation. We generally denote states by the symbol i and successor states by the symbol j . At state i , the use of a control u specifies the transition probability $p_{ij}(u)$ to the next state j . The state i is an element of a finite state space, and the control u is constrained to take values in a given finite constraint set $U(i)$, which may depend on the current state i .

As discussed in Section 1.1, the underlying system equation is

$$x_{k+1} = w_k,$$

where w_k is the disturbance. We will generally suppress w_k from the cost to simplify notation. Thus we will assume a k th stage cost $g(x_k, u_k)$ for using control u_k at state x_k . This amounts to averaging the cost per stage over all successor states, which makes no essential difference in the subsequent analysis. Thus, if $\tilde{g}(i, u, j)$ is the cost of using u at state i and moving to state j , we use as cost per stage the expected cost $g(i, u)$ given by

$$g(i, u) = \sum_j p_{ij}(u) \tilde{g}(i, u, j).$$

The total expected cost associated with an initial state i and an admissible policy $\pi = \{\mu_0, \mu_1, \dots\}$ [one with $\mu_k(i) \in U(i)$ for all i and k] is

$$J_\pi(i) = \lim_{N \rightarrow \infty} E \left\{ \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k)) \mid x_0 = i \right\},$$

where α is a discount factor with $0 < \alpha \leq 1$. In the following two sections, we will impose assumptions that guarantee the existence of the above limit. The optimal cost from state i , i.e., the minimum of $J_\pi(i)$ over all admissible π , is denoted by $J^*(i)$. A *stationary policy* is an admissible policy of the form $\pi = \{\mu, \mu, \dots\}$, and its corresponding cost function is denoted by $J_\mu(i)$. For brevity, we refer to $\{\mu, \mu, \dots\}$ as the stationary policy μ . We say that μ is optimal if

$$J_\mu(i) = J^*(i) = \min_{\pi} J_\pi(i), \quad \text{for all states } i.$$

5.2 STOCHASTIC SHORTEST PATH PROBLEMS

Here, we assume that there is no discounting ($\alpha = 1$), and that *there is a special cost-free termination state t* . Once the system reaches that state, it remains there at no further cost, i.e., $p_{tt}(u) = 1$ and $g(t, u) = 0$ for all $u \in U(t)$. We denote by $1, \dots, n$ the states other than the termination state t ; see Fig. 5.2.1.

We are interested in problems where reaching the termination state is inevitable, at least under an optimal policy. Thus, the essence of the

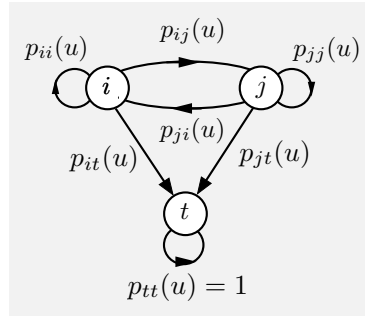


Figure 5.2.1 The transition graph of a stochastic shortest path problem. There are n states, plus the termination state t , with transition probabilities $p_{ij}(u)$. The destination is cost-free and absorbing.

problem is to reach the termination state with minimum expected cost. We call this the *stochastic shortest path problem*. The deterministic shortest path problem is obtained as the special case where for each state-control pair (i, u) , the transition probability $p_{ij}(u)$ is equal to 1 for a unique state j that depends on (i, u) . The reader may also verify that the finite horizon problem of Chapter 1 can be obtained as a special case by viewing as state the pair (x_k, k) and lumping all pairs (x_N, N) into a termination state t .

Certain conditions are required to guarantee that, at least under an optimal policy, termination occurs with probability 1. We will make the following assumption that guarantees eventual termination under *all* policies:

Assumption 5.2.1: There exists an integer m such that regardless of the policy used and the initial state, there is positive probability that the termination state will be reached after no more than m stages; i.e., for all admissible policies π we have

$$\rho_\pi = \max_{i=1, \dots, n} P\{x_m \neq t \mid x_0 = i, \pi\} < 1. \quad (5.3)$$

We note, however, that the results to be presented are valid under more general circumstances.[†] Furthermore, it can be shown that if there

[†] Extensions of the stochastic shortest path theory are given in Vol. II. For a brief summary, let us call a stationary policy π *proper* if the condition (5.3) is satisfied for some m , and call π *improper* otherwise. It can be shown that Assumption 5.2.1 is equivalent to the seemingly weaker assumption that all stationary policies are proper (see Vol. II, Exercise 3.6). However, the results of the

exists an integer m with the property of Assumption 5.2.1, then there also exists an integer less or equal to n with this property (Exercise 5.12). Thus, we can always use $m = n$ in Assumption 5.2.1, if no smaller value of m is known. Let

$$\rho = \max_{\pi} \rho_{\pi}.$$

Note that ρ_{π} depends only on the first m components of the policy π . Furthermore, since the number of controls available at each state is finite, the number of distinct m -stage policies is also finite. It follows that there can be only a finite number of distinct values of ρ_{π} so that

$$\rho < 1.$$

We therefore have for any π and any initial state i

$$\begin{aligned} P\{x_{2m} \neq t \mid x_0 = i, \pi\} &= P\{x_{2m} \neq t \mid x_m \neq t, x_0 = i, \pi\} \\ &\quad \cdot P\{x_m \neq t \mid x_0 = i, \pi\} \\ &\leq \rho^2. \end{aligned}$$

More generally, for each admissible policy π , the probability of not reaching the termination state after km stages diminishes like ρ^k regardless of the initial state, i.e.,

$$P\{x_{km} \neq t \mid x_0 = i, \pi\} \leq \rho^k, \quad i = 1, \dots, n. \quad (5.4)$$

Thus the limit defining the associated total cost vector J_{π} exists and is finite, since the expected cost incurred in the m periods between km and

subsequent Prop. 5.2.1 can also be shown under the genuinely weaker assumption that there exists at least one proper policy, and furthermore, every improper policy results in infinite expected cost from at least one initial state (see [BeT89], [BeT91], or Vol. II, Chapter 3). These assumptions, when specialized to deterministic shortest path problems, are similar to the assumptions of Chapter 2. They imply that there is at least one path to the destination from every starting node and that all cycles have positive cost.

A weaker set of assumptions, whereby improper policies are allowed to have finite expected cost and even be optimal [but the optimal costs $J^*(i)$ are assumed finite for all i], is considered in the papers [Ber15b] and [BeY16]. While some serious complications may occur under these weaker assumptions (including issues of validity of Bellman's equation, and various algorithms; see Exercise 5.29 for an example), the results of Prop. 5.2.1 can be shown in a weaker form.

Still another set of assumptions under which the results of Prop. 5.2.1 hold is described in Exercise 5.28, where again improper policies are allowed, but the stage costs $g(i, u)$ are assumed nonnegative, and the optimal costs $J^*(i)$ are assumed finite. Finally, we note that there is a parallel theory for minimax-type shortest path problems (also called *robust shortest path problems*), which involve a set-membership description of uncertainty; see [Ber14].

$(k+1)m-1$ is bounded in absolute value by

$$m\rho^k \max_{\substack{i=1,\dots,n \\ u \in U(i)}} |g(i, u)|.$$

In particular, we have

$$|J_\pi(i)| \leq \sum_{k=0}^{\infty} m\rho^k \max_{\substack{i=1,\dots,n \\ u \in U(i)}} |g(i, u)| = \frac{m}{1-\rho} \max_{\substack{i=1,\dots,n \\ u \in U(i)}} |g(i, u)|. \quad (5.5)$$

The results of the following proposition are basic and are typical of many infinite horizon problems. Despite its length (and the lengthy expressions it involves), the proof is not complicated. The key idea is that the “tail” of the cost series,

$$\sum_{k=mK}^{\infty} E\{g(x_k, \mu_k(x_k))\},$$

vanishes as K increases to ∞ , since the probability that $x_{mK} \neq t$ decreases like ρ^K [cf. Eq. (5.4)].

Proposition 5.2.1: Under Assumption 5.2.1, the following hold for the stochastic shortest path problem:

- (a) Given any initial conditions $J_0(1), \dots, J_0(n)$, the sequence $J_k(i)$ generated by the iteration

$$J_{k+1}(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_k(j) \right], \quad i = 1, \dots, n, \quad (5.6)$$

converges to the optimal cost $J^*(i)$ for each i . [Note that, by reversing the time index this iteration can be viewed as the DP algorithm for a finite horizon problem with terminal cost function equal to J_0 . In fact, $J_k(i)$ is the optimal cost starting from state i of a k -stage problem with cost per stage given by g and terminal cost at the end of the k stages given by J_0 .]

- (b) The optimal costs $J^*(1), \dots, J^*(n)$ satisfy Bellman’s equation,

$$J^*(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J^*(j) \right], \quad i = 1, \dots, n, \quad (5.7)$$

and in fact they are the unique solution of this equation.

- (c) For any stationary policy μ , the costs $J_\mu(1), \dots, J_\mu(n)$ are the unique solution of the equation

$$J_\mu(i) = g(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) J_\mu(j), \quad i = 1, \dots, n.$$

Furthermore, given any initial conditions $J_0(1), \dots, J_0(n)$, the sequence $J_k(i)$ generated by the DP iteration

$$J_{k+1}(i) = g(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) J_k(j), \quad i = 1, \dots, n,$$

converges to the cost $J_\mu(i)$ for each i .

- (d) A stationary policy μ is optimal if and only if for every state i , $\mu(i)$ attains the minimum in Bellman's equation (5.7).

Proof: (a) For every positive integer K , initial state x_0 , and policy $\pi = \{\mu_0, \mu_1, \dots\}$, we break down the cost $J_\pi(x_0)$ into the portions incurred over the first mK stages and over the remaining stages:

$$\begin{aligned} J_\pi(x_0) &= \lim_{N \rightarrow \infty} E \left\{ \sum_{k=0}^{N-1} g(x_k, \mu_k(x_k)) \right\} \\ &= E \left\{ \sum_{k=0}^{mK-1} g(x_k, \mu_k(x_k)) \right\} + \lim_{N \rightarrow \infty} E \left\{ \sum_{k=mK}^{N-1} g(x_k, \mu_k(x_k)) \right\}. \end{aligned}$$

Let B denote the following upper bound on the cost of an m -stage cycle, assuming termination does not occur during the cycle,

$$B = m \max_{\substack{i=1, \dots, n \\ u \in U(i)}} |g(i, u)|.$$

The expected cost during the K th m -stage cycle [stages Km to $(K+1)m-1$] is upper bounded by $B\rho^K$ [cf. Eqs. (5.4) and (5.5)], so that

$$\left| \lim_{N \rightarrow \infty} E \left\{ \sum_{k=mK}^{N-1} g(x_k, \mu_k(x_k)) \right\} \right| \leq B \sum_{k=K}^{\infty} \rho^k = \frac{\rho^K B}{1 - \rho}.$$

Also, denoting $J_0(t) = 0$, let us view J_0 as a terminal cost function and bound its expected value under π after mK stages. We have

$$\begin{aligned} |E\{J_0(x_{mK})\}| &= \left| \sum_{i=1}^n P(x_{mK} = i \mid x_0, \pi) J_0(i) \right| \\ &\leq \left(\sum_{i=1}^n P(x_{mK} = i \mid x_0, \pi) \right) \max_{i=1, \dots, n} |J_0(i)| \\ &\leq \rho^K \max_{i=1, \dots, n} |J_0(i)|, \end{aligned}$$

since the probability that $x_{mK} \neq t$ is less or equal to ρ^K for any policy. Combining the preceding relations, we obtain

$$\begin{aligned} & -\rho^K \max_{i=1,\dots,n} |J_0(i)| + J_\pi(x_0) - \frac{\rho^K B}{1-\rho} \\ & \leq E \left\{ J_0(x_{mK}) + \sum_{k=0}^{mK-1} g(x_k, \mu_k(x_k)) \right\} \\ & \leq \rho^K \max_{i=1,\dots,n} |J_0(i)| + J_\pi(x_0) + \frac{\rho^K B}{1-\rho}. \end{aligned} \quad (5.8)$$

Note that the expected value in the middle term of the above inequalities is the mK -stage cost of policy π starting from state x_0 , with a terminal cost $J_0(x_{mK})$; the minimum of this cost over all π is equal to the value $J_{mK}(x_0)$, which is generated by the DP recursion (5.6) after mK iterations. Thus, by taking the minimum over π in Eq. (5.8), we obtain for all x_0 and K ,

$$\begin{aligned} & -\rho^K \max_{i=1,\dots,n} |J_0(i)| + J^*(x_0) - \frac{\rho^K B}{1-\rho} \\ & \leq J_{mK}(x_0) \\ & \leq \rho^K \max_{i=1,\dots,n} |J_0(i)| + J^*(x_0) + \frac{\rho^K B}{1-\rho}, \end{aligned} \quad (5.9)$$

and by taking the limit as $K \rightarrow \infty$, we obtain

$$\lim_{K \rightarrow \infty} J_{mK}(x_0) = J^*(x_0)$$

for all x_0 . Since

$$|J_{mK+\ell}(x_0) - J_{mK}(x_0)| \leq \rho^K B, \quad \ell = 1, \dots, m,$$

we see that $\lim_{K \rightarrow \infty} J_{mK+\ell}(x_0)$ is the same for all $\ell = 1, \dots, m$, so that

$$\lim_{k \rightarrow \infty} J_k(x_0) = J^*(x_0).$$

(b) By taking the limit as $k \rightarrow \infty$ in the DP iteration (5.6) and using the result of part (a), we see that $J^*(1), \dots, J^*(n)$ satisfy Bellman's equation (we are using here the fact that the limit and minimization operations commute when the minimization is over a finite number of alternatives). To show uniqueness, observe that if $J(1), \dots, J(n)$ satisfy Bellman's equation, then the DP iteration (5.6) starting from $J(1), \dots, J(n)$ just replicates $J(1), \dots, J(n)$. It follows from the convergence result of part (a) that $J(i) = J^*(i)$ for all i .

(c) Given the stationary policy μ , we can consider a modified stochastic shortest path problem, which is the same as the original except that the

control constraint set contains only one element for each state i , the control $\mu(i)$; i.e., the control constraint set is $\tilde{U}(i) = \{\mu(i)\}$ instead of $U(i)$. From part (b) we then obtain that $J_\mu(1), \dots, J_\mu(n)$ solve uniquely Bellman's equation for this modified problem, i.e.,

$$J_\mu(i) = g(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) J_\mu(j), \quad i = 1, \dots, n.$$

From part (a), the corresponding DP iteration converges to $J_\mu(i)$.

(d) We have that $\mu(i)$ attains the minimum in Eq. (5.7) if and only if

$$\begin{aligned} J^*(i) &= \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J^*(j) \right] \\ &= g(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) J^*(j), \quad i = 1, \dots, n. \end{aligned}$$

Part (c) and this equation imply that $J_\mu(i) = J^*(i)$ for all i . Conversely, if $J_\mu(i) = J^*(i)$ for all i , parts (b) and (c) imply this equation. **Q.E.D.**

Example 5.2.1 (Minimizing Expected Time to Termination)

The case where

$$g(i, u) = 1, \quad i = 1, \dots, n, \quad u \in U(i),$$

corresponds to a problem where the objective is to terminate as fast as possible on the average, while the corresponding optimal cost $J^*(i)$ is the minimum expected time to termination starting from state i . Under our assumptions, the costs $J^*(i)$ uniquely solve Bellman's equation, which has the form

$$J^*(i) = \min_{u \in U(i)} \left[1 + \sum_{j=1}^n p_{ij}(u) J^*(j) \right], \quad i = 1, \dots, n.$$

In the special case where there is only one control at each state, $J^*(i)$ represents the mean first passage time from i to t (see Appendix D). These times, denoted m_i , are the unique solution of the equations

$$m_i = 1 + \sum_{j=1}^n p_{ij} m_j, \quad i = 1, \dots, n.$$

Example 5.2.2

A spider and a fly move along a straight line at times $k = 0, 1, \dots$. The initial positions of the fly and the spider are integer. At each time period, the fly moves one unit to the left with probability p , one unit to the right with probability p , and stays where it is with probability $1 - 2p$. The spider, knows the position of the fly at the beginning of each period, and will always move one unit towards the fly if its distance from the fly is more than one unit. If the spider is one unit away from the fly, it will either move one unit towards the fly or stay where it is. If the spider and the fly land in the same position at the end of a period, then the spider captures the fly and the process terminates. The spider's objective is to capture the fly in minimum expected time.

We view as state the distance between spider and fly. Then the problem can be formulated as a stochastic shortest path problem with states $0, 1, \dots, n$, where n is the initial distance. State 0 is the termination state where the spider captures the fly. Let us denote $p_{1j}(M)$ and $p_{1j}(\overline{M})$ the transition probabilities from state 1 to state j if the spider moves and does not move, respectively, and let us denote by p_{ij} the transition probabilities from a state $i \geq 2$. We have

$$p_{ii} = p, \quad p_{i(i-1)} = 1 - 2p, \quad p_{i(i-2)} = p, \quad i \geq 2,$$

$$p_{11}(M) = 2p, \quad p_{10}(M) = 1 - 2p,$$

$$p_{12}(\overline{M}) = p, \quad p_{11}(\overline{M}) = 1 - 2p, \quad p_{10}(\overline{M}) = p,$$

with all other transition probabilities being 0.

For states $i \geq 2$, Bellman's equation is written as

$$J^*(i) = 1 + pJ^*(i) + (1 - 2p)J^*(i - 1) + pJ^*(i - 2), \quad i \geq 2, \quad (5.10)$$

where $J^*(0) = 0$ by definition. The only state where the spider has a choice is when it is one unit away from the fly, and for that state Bellman's equation is given by

$$J^*(1) = 1 + \min[2pJ^*(1), pJ^*(2) + (1 - 2p)J^*(1)], \quad (5.11)$$

where the first and the second expression within the bracket above are associated with the spider moving and not moving, respectively. By writing Eq. (5.10) for $i = 2$, we obtain

$$J^*(2) = 1 + pJ^*(2) + (1 - 2p)J^*(1),$$

from which

$$J^*(2) = \frac{1}{1 - p} + \frac{(1 - 2p)J^*(1)}{1 - p}. \quad (5.12)$$

Substituting this expression in Eq. (5.11), we obtain

$$J^*(1) = 1 + \min \left[2pJ^*(1), \frac{p}{1 - p} + \frac{p(1 - 2p)J^*(1)}{1 - p} + (1 - 2p)J^*(1) \right],$$

or equivalently,

$$J^*(1) = 1 + \min \left[2pJ^*(1), \frac{p}{1-p} + \frac{(1-2p)J^*(1)}{1-p} \right].$$

To solve the above equation, we consider the two cases where the first expression within the bracket is larger and is smaller than the second expression. Thus we solve for $J^*(1)$ in the two cases where

$$J^*(1) = 1 + 2pJ^*(1), \quad (5.13)$$

$$2pJ^*(1) \leq \frac{p}{1-p} + \frac{(1-2p)J^*(1)}{1-p}, \quad (5.14)$$

and

$$J^*(1) = 1 + \frac{p}{1-p} + \frac{(1-2p)J^*(1)}{1-p}, \quad (5.15)$$

$$2pJ^*(1) \geq \frac{p}{1-p} + \frac{(1-2p)J^*(1)}{1-p}. \quad (5.16)$$

The solution of Eq. (5.13) is seen to be

$$J^*(1) = \frac{1}{1-2p},$$

and by substitution in Eq. (5.14), we find that this solution is valid when

$$\frac{2p}{1-2p} \leq \frac{p}{1-p} + \frac{1}{1-p},$$

or equivalently (after some calculation), $p \leq 1/3$. Thus for $p \leq 1/3$, it is optimal for the spider to move when it is one unit away from the fly.

Similarly, the solution of Eq. (5.15) is seen to be

$$J^*(1) = \frac{1}{p},$$

and by substitution in Eq. (5.16), we find that this solution is valid when

$$2 \geq \frac{p}{1-p} + \frac{1-2p}{p(1-p)},$$

or equivalently (after some calculation), $p \geq 1/3$. Thus, for $p \geq 1/3$ it is optimal for the spider not to move when it is one unit away from the fly.

The minimal expected number of steps for capture when the spider is one unit away from the fly was calculated earlier to be

$$J^*(1) = \begin{cases} 1/(1-2p) & \text{if } p \leq 1/3, \\ 1/p & \text{if } p \geq 1/3. \end{cases}$$

Given the value of $J^*(1)$, we can calculate from Eq. (5.12) the minimal expected number of steps for capture when two units away, $J^*(2)$, and we can then obtain the remaining values $J^*(i)$, $i = 3, \dots, n$, from Eq. (5.10).

5.3 COMPUTATIONAL METHODS

We now turn to computational methods for stochastic shortest path problems. There are three major types of methods, which are described in the next three subsections.

5.3.1 Value Iteration

The DP iteration

$$J_{k+1}(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_k(j) \right], \quad i = 1, \dots, n, \quad (5.17)$$

is called *value iteration* and is a principal method for calculating the optimal cost function J^* . Generally, value iteration requires an infinite number of iterations, although there are important special cases of stochastic shortest path problems where it terminates finitely (see Vol. II, Section 3.4). Note that from Eq. (5.9) we obtain that the error

$$|J_{mK}(i) - J^*(i)|$$

is bounded by a constant multiple of ρ^K .

The value iteration algorithm can sometimes be strengthened with the use of some error bounds. In particular, it can be shown (see Exercise 5.13) that for all k and j , we have

$$J_{k+1}(j) + (N^*(j) - 1) \underline{c}_k \leq J^*(j) \leq J_{\mu^k}(j) \leq J_{k+1}(j) + (N^k(j) - 1) \bar{c}_k, \quad (5.18)$$

where the policy μ^k is such that $\mu^k(i)$ attains the minimum in the k th value iteration (5.17) for all i , and

$N^*(j)$: The average number of stages to reach t starting from j and using some optimal stationary policy,

$N^k(j)$: The average number of stages to reach t starting from j and using the stationary policy μ^k ,

$$\underline{c}_k = \min_{i=1, \dots, n} [J_{k+1}(i) - J_k(i)], \quad \bar{c}_k = \max_{i=1, \dots, n} [J_{k+1}(i) - J_k(i)].$$

Unfortunately, the values $N^*(j)$ and $N^k(j)$ are easily computed or approximated only in the presence of special problem structure. Despite this fact, the bounds (5.18) often provide a useful guideline for stopping the value iteration algorithm while being assured that J_k approximates J^* with sufficient accuracy. Moreover, under suitable conditions, it may be possible to provide bounds to the left and right sides of Eq. (5.18), which are more easily computable, as shown in the related works [HaA15] and [Han17].

5.3.2 Policy Iteration

An alternative to value iteration is *policy iteration*. This algorithm starts with a stationary policy μ^0 , and generates iteratively a sequence of new policies μ^1, μ^2, \dots as follows:

Given the policy μ^k , we perform a *policy evaluation step*, that computes $J_{\mu^k}(i)$, $i = 1, \dots, n$, as the solution of the (linear) system of equations

$$J(i) = g(i, \mu^k(i)) + \sum_{j=1}^n p_{ij}(\mu^k(i)) J(j), \quad i = 1, \dots, n, \quad (5.19)$$

in the n unknowns $J(1), \dots, J(n)$ [cf. Prop. 5.2.1(c)]. We then perform a *policy improvement step*, which computes a new policy μ^{k+1} as

$$\mu^{k+1}(i) = \arg \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_{\mu^k}(j) \right], \quad i = 1, \dots, n. \quad (5.20)$$

The process is repeated with μ^{k+1} used in place of μ^k , unless we have $J_{\mu^{k+1}}(i) = J_{\mu^k}(i)$ for all i , in which case the algorithm terminates with the policy μ^k . The following proposition establishes the validity of policy iteration, including finite termination with an optimal policy.

Proposition 5.3.1: Under Assumption 5.2.1, the policy iteration algorithm for the stochastic shortest path problem generates an improving sequence of policies [i.e., $J_{\mu^{k+1}}(i) \leq J_{\mu^k}(i)$ for all i and k] and terminates with an optimal policy.

Proof: For any k , consider the sequence generated by the recursion

$$J_{N+1}(i) = g(i, \mu^{k+1}(i)) + \sum_{j=1}^n p_{ij}(\mu^{k+1}(i)) J_N(j), \quad i = 1, \dots, n,$$

where $N = 0, 1, \dots$, and

$$J_0(i) = J_{\mu^k}(i), \quad i = 1, \dots, n.$$

From Eqs. (5.19) and (5.20), we have

$$\begin{aligned} J_0(i) &= g(i, \mu^k(i)) + \sum_{j=1}^n p_{ij}(\mu^k(i)) J_0(j) \\ &\geq g(i, \mu^{k+1}(i)) + \sum_{j=1}^n p_{ij}(\mu^{k+1}(i)) J_0(j) \\ &= J_1(i), \end{aligned}$$

for all i . By using the above inequality we obtain (compare with the monotonicity property of DP, Exercise 1.28 in Chapter 1)

$$\begin{aligned} J_1(i) &= g(i, \mu^{k+1}(i)) + \sum_{j=1}^n p_{ij}(\mu^{k+1}(i)) J_0(j) \\ &\geq g(i, \mu^{k+1}(i)) + \sum_{j=1}^n p_{ij}(\mu^{k+1}(i)) J_1(j) \\ &= J_2(i), \end{aligned}$$

for all i , and by continuing similarly we have

$$J_0(i) \geq J_1(i) \geq \cdots \geq J_N(i) \geq J_{N+1}(i) \geq \cdots, \quad i = 1, \dots, n. \quad (5.21)$$

Since by Prop. 5.2.1(c), $J_N(i) \rightarrow J_{\mu^{k+1}}(i)$, we obtain $J_0(i) \geq J_{\mu^{k+1}}(i)$ or

$$J_{\mu^k}(i) \geq J_{\mu^{k+1}}(i), \quad i = 1, \dots, n, \quad k = 0, 1, \dots$$

Thus the sequence of generated policies is improving, and since the number of stationary policies is finite, we must after a finite number of iterations, say $k + 1$, obtain $J_{\mu^k}(i) = J_{\mu^{k+1}}(i)$ for all i . Then we will have equality throughout in Eq. (5.21), which means that

$$J_{\mu^k}(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_{\mu^k}(j) \right], \quad i = 1, \dots, n.$$

Thus the costs $J_{\mu^k}(1), \dots, J_{\mu^k}(n)$ solve Bellman's equation, and by Prop. 5.2.1(b), it follows that $J_{\mu^k}(i) = J^*(i)$ and that μ^k is optimal. **Q.E.D.**

The linear system of equations (5.19) of the policy evaluation step can be solved by standard methods such as Gaussian elimination, but when the number of states is large, this is cumbersome and time-consuming. A typically more efficient alternative is to approximate the policy evaluation step with a few value iterations aimed at solving the corresponding system (5.19). One can show that the policy iteration method that uses such approximate policy evaluation yields in the limit the optimal costs and an optimal stationary policy, even if we evaluate each policy using an arbitrary positive number of value iterations. This variant of the policy iteration method is called *optimistic policy iteration* (or sometimes *modified policy iteration*), and is discussed in more detail in Vol. II, starting with Section 2.3.

Another possibility for approximating the policy evaluation step is to use simulation, and this is a key idea in the rollout algorithm, to be discussed in Section 6.4. Simulation also plays an important role in the

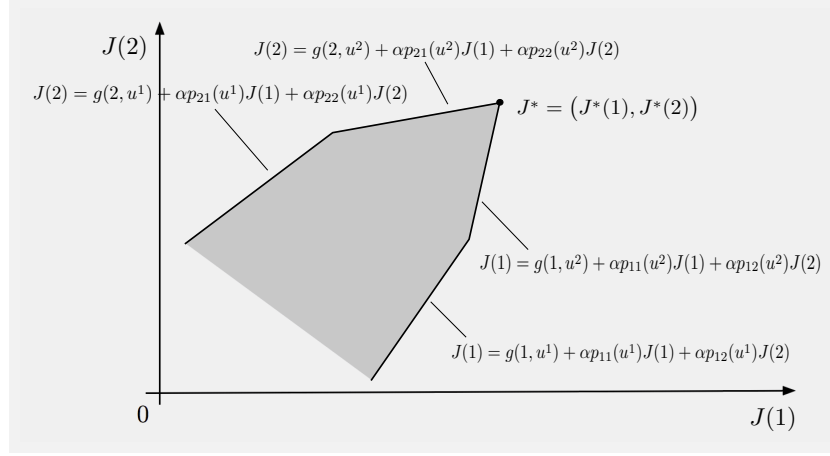


Figure 5.3.1 Linear program associated with a two-state stochastic shortest path problem. The constraint set is shaded, and the objective to maximize is $J(1) + J(2)$. Note that because we have $J(i) \leq J^*(i)$ for all i and vectors J in the constraint set, the vector J^* maximizes any linear cost function of the form $\sum_{i=1}^n \beta_i J(i)$ where $\beta_i \geq 0$ for all i . If $\beta_i > 0$ for all i , then J^* is the unique optimal solution of the corresponding linear program.

approximate DP methodology, discussed in Chapters 6 and 7 of Vol. II. In particular, when the number of states is large, one can try to approximate the cost-to-go function J_{μ^k} by simulating a number of trajectories under the policy μ^k , and perform some form of least squares fit of J_{μ^k} using an approximation architecture (see Section 6.3). These are a number of variations of this idea, involving different types of simulation, a variety of policy iteration-type algorithms, and optimistic versions thereof, which are discussed in more detail in Vol. II, and in the research monograph [BeT96].

5.3.3 Linear Programming

Suppose that we use value iteration to generate a sequence of vectors $J_k = (J_k(1), \dots, J_k(n))$ starting with an initial condition vector $J_0 = (J_0(1), \dots, J_0(n))$ such that

$$J_0(i) \leq \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_0(j) \right], \quad i = 1, \dots, n.$$

Then we will have $J_k(i) \leq J_{k+1}(i)$ for all k and i (the monotonicity property of DP; Exercise 1.28 in Chapter 1). It follows from Prop. 5.2.1(a) that we will also have $J_0(i) \leq J^*(i)$ for all i . Thus J^* is the “largest” J that

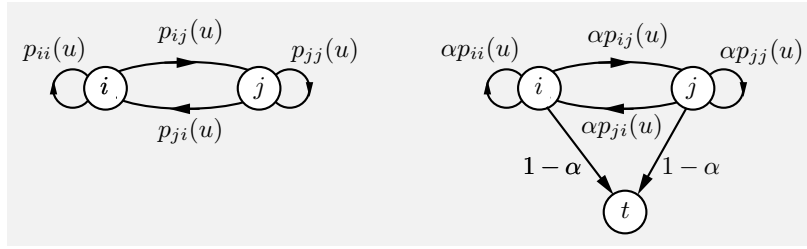


Figure 5.4.1 Transition probabilities for an α -discounted problem and its associated stochastic shortest path problem. In the latter problem, the probability that the state is not t after k stages is α^k . The expected cost at each state $i = 1, \dots, n$ is $g(i, u)$ for both problems, but it must be multiplied by α^k because of discounting (in the discounted case) or because it is incurred with probability α^k when termination has not yet been reached (in the stochastic shortest path case).

satisfies the constraint

$$J(i) \leq g(i, u) + \sum_{j=1}^n p_{ij}(u) J(j), \quad \text{for all } i = 1, \dots, n \text{ and } u \in U(i). \quad (5.22)$$

In particular, $J^*(1), \dots, J^*(n)$ solve the linear program of maximizing $\sum_{i=1}^n J(i)$ subject to the constraint (5.22) (see Fig. 5.3.1). Unfortunately, for large n the dimension of this program can be very large and its solution can be impractical, particularly in the absence of special structure. We refer to Section 2.4 of Vol. II for further discussion.

5.4 DISCOUNTED PROBLEMS

We now consider a discounted problem, where there is a discount factor $\alpha < 1$. We will show that this problem can be converted to a stochastic shortest path problem for which the analysis of the preceding section holds. To see this, let $i = 1, \dots, n$ be the states, and consider an associated stochastic shortest path problem involving the states $1, \dots, n$ plus an extra termination state t , with state transitions and costs obtained as follows: From a state $i \neq t$, when control u is applied, a cost $g(i, u)$ is incurred, and the next state is j with probability $\alpha p_{ij}(u)$ and t with probability $1 - \alpha$; see Fig. 5.4.1. Note that Assumption 5.2.1 of the preceding section is satisfied for this stochastic shortest path problem.

Suppose now that we use the same policy in the discounted problem and in the associated stochastic shortest path problem. Then, as long as termination has not occurred, the state evolution in the two problems is governed by the same transition probabilities. Furthermore, the expected cost of the k th stage of the associated shortest path problem is

$g(x_k, \mu_k(x_k))$ multiplied by the probability that state t has not yet been reached, which is α^k . This is also the expected cost of the k th stage for the discounted problem. Thus the cost of any policy starting from a given state, is the same for the original discounted problem and for the associated stochastic shortest path problem. Furthermore, value iteration produces identical iterates for the two problems. We can thus apply the results of the preceding section to the latter problem and obtain the following:

Proposition 5.4.1: The following hold for the discounted problem:

- (a) The value iteration algorithm

$$J_{k+1}(i) = \min_{u \in U(i)} \left[g(i, u) + \alpha \sum_{j=1}^n p_{ij}(u) J_k(j) \right], \quad i = 1, \dots, n, \quad (5.23)$$

converges to the optimal costs $J^*(i)$, $i = 1, \dots, n$, starting from arbitrary initial conditions $J_0(1), \dots, J_0(n)$.

- (b) The optimal costs $J^*(1), \dots, J^*(n)$ of the discounted problem satisfy Bellman's equation,

$$J^*(i) = \min_{u \in U(i)} \left[g(i, u) + \alpha \sum_{j=1}^n p_{ij}(u) J^*(j) \right], \quad i = 1, \dots, n, \quad (5.24)$$

and in fact they are the unique solution of this equation.

- (c) For any stationary policy μ , the costs $J_\mu(1), \dots, J_\mu(n)$ are the unique solution of the equation

$$J_\mu(i) = g(i, \mu(i)) + \alpha \sum_{j=1}^n p_{ij}(\mu(i)) J_\mu(j), \quad i = 1, \dots, n.$$

Furthermore, given any initial conditions $J_0(1), \dots, J_0(n)$, the sequence $J_k(i)$ generated by the DP iteration

$$J_{k+1}(i) = g(i, \mu(i)) + \alpha \sum_{j=1}^n p_{ij}(\mu(i)) J_k(j), \quad i = 1, \dots, n,$$

converges to the cost $J_\mu(i)$ for each i .

- (d) A stationary policy μ is optimal if and only if for every state i , $\mu(i)$ attains the minimum in Bellman's equation (5.24).

(e) The policy iteration algorithm given by

$$\mu^{k+1}(i) = \arg \min_{u \in U(i)} \left[g(i, u) + \alpha \sum_{j=1}^n p_{ij}(u) J_{\mu^k}(j) \right], \quad i = 1, \dots, n,$$

generates an improving sequence of policies and terminates with an optimal policy.

Proof: Parts (a)-(d) and part (e) are proved by applying parts (a)-(d) of Prop. 5.2.1, and Prop. 5.3.1, respectively, to the associated stochastic shortest path problem described above. **Q.E.D.**

Bellman's equation (5.24) has a familiar DP interpretation. At state i , the optimal cost $J^*(i)$ is the minimum over all controls of the sum of the expected current stage cost and the expected optimal cost of all future stages. The former cost is $g(i, u)$. The latter cost is $J^*(j)$, but since this cost starts accumulating after one stage, it is discounted by multiplication with α .

As in the case of stochastic shortest path problems [see Eq. (5.9) and the discussion following the proof of Prop. 5.2.1], we can show that the error

$$|J_k(i) - J^*(i)|$$

is bounded by a constant times α^k . Furthermore, the error bounds (5.18) become

$$J_{k+1}(j) + \frac{\alpha}{1-\alpha} \underline{c}_k \leq J^*(j) \leq J_{\mu^k}(j) \leq J_{k+1}(j) + \frac{\alpha}{1-\alpha} \bar{c}_k, \quad (5.25)$$

where μ^k is such that $\mu^k(i)$ attains the minimum in the k th value iteration (5.23) for all i , and

$$\underline{c}_k = \min_{i=1, \dots, n} [J_{k+1}(i) - J_k(i)], \quad \bar{c}_k = \max_{i=1, \dots, n} [J_{k+1}(i) - J_k(i)],$$

since for the associated stochastic shortest path problem it can be shown that for every policy and starting state, the expected number of stages to reach the termination state t is $1/(1-\alpha)$, so that the terms $N^*(j) - 1$ and $N^k(j) - 1$ appearing in Eq. (5.18) are equal to $\alpha/(1-\alpha)$. We note also that there are a number of additional enhancements to the value iteration algorithm for the discounted problem (see Section 2.2 of Vol. II). There are also discounted cost variants of the optimistic policy iteration and linear programming approaches discussed for stochastic shortest path problems.

Example 5.4.1 (Asset Selling)

Consider an infinite horizon version of the asset selling example of Section 3.4, assuming the set of possible offers is finite. Here, if accepted, the amount x_k offered in period k , will be invested at a rate of interest r . By depreciating the sale amount to period 0 dollars, we view $(1+r)^{-k}x_k$ as the reward for selling the asset in period k at a price x_k , where $r > 0$ is the rate of interest. Then we have a total discounted reward problem with discount factor $\alpha = 1/(1+r)$. The analysis of the present section is applicable, and the optimal value function J^* is the unique solution of Bellman's equation

$$J^*(x) = \max \left[x, \frac{E\{J^*(w)\}}{1+r} \right],$$

(see Section 3.4). The optimal reward function is characterized by the critical number

$$\bar{\alpha} = \frac{E\{J^*(w)\}}{1+r},$$

which can be calculated as in Section 3.4. An optimal policy is to sell if and only if the current offer x_k is greater than or equal to $\bar{\alpha}$.

Example 5.4.2

A manufacturer at each time period receives an order for her product with probability p and receives no order with probability $1-p$. At any period she has a choice of processing all the unfilled orders in a batch, or process no order at all. The cost per unfilled order at each time period is $c > 0$, and the setup cost to process the unfilled orders is $K > 0$. The manufacturer wants to find a processing policy that minimizes the total expected cost, assuming the discount factor is $\alpha < 1$ and the maximum number of orders that can remain unfilled is n .

Here the state is the number of unfilled orders at the beginning of each period, and Bellman's equation takes the form

$$J^*(i) = \min [K + \alpha(1-p)J^*(0) + \alpha p J^*(1), ci + \alpha(1-p)J^*(i) + \alpha p J^*(i+1)], \quad (5.26)$$

for the states $i = 0, 1, \dots, n-1$, and takes the form

$$J^*(n) = K + \alpha(1-p)J^*(0) + \alpha p J^*(1) \quad (5.27)$$

for state n . The first expression within brackets in Eq. (5.26) corresponds to processing the i unfilled orders, while the second expression corresponds to leaving the orders unfilled for one more period. When the maximum n of unfilled orders is reached, the orders must necessarily be processed, as indicated by Eq. (5.27).

To solve the problem, we observe that the optimal cost $J^*(i)$ is monotonically nondecreasing in i . This is intuitively clear, and can be rigorously

proved by using the value iteration method. In particular, we can show by using the (finite horizon) DP algorithm that the k -stage optimal cost functions $J_k(i)$ are monotonically nondecreasing in i for all k (Exercise 5.7), and then argue that the optimal infinite horizon cost function $J^*(i)$ is also monotonically nondecreasing in i , since

$$J^*(i) = \lim_{k \rightarrow \infty} J_k(i)$$

by Prop. 5.4.1(a). Given that $J^*(i)$ is monotonically nondecreasing in i , from Eq. (5.26) we have that if processing a batch of m orders is optimal, i.e.,

$$K + \alpha(1 - p)J^*(0) + \alpha p J^*(1) \leq cm + \alpha(1 - p)J^*(m) + \alpha p J^*(m + 1),$$

then processing a batch of $m + 1$ orders is also optimal. Therefore a *threshold policy*, i.e., a policy that processes the orders if their number exceeds some threshold integer m^* , is optimal.

We leave it as Exercise 5.8 for the reader to verify that if we start the policy iteration algorithm with a threshold policy, every subsequently generated policy will be a threshold policy. Since there are $n + 1$ distinct threshold policies, and the sequence of generated policies is improving, it follows that the policy iteration algorithm will yield an optimal policy after at most n iterations.

5.5 AVERAGE COST PER STAGE PROBLEMS

5.6 SEMI-MARKOV PROBLEMS

5.7 NOTES, SOURCES, AND EXERCISES

In this chapter we have provided a gentle introduction to infinite horizon problems. There is an extensive theory for these problems with interesting and challenging mathematical and computational content. Volume II provides a comprehensive treatment and gives many references to the literature. Moreover, Vol. II provides an extensive account of approximate solution methods that are suitable for large-scale DP problems, where the value and policy iteration algorithms cannot be used in the form given here, because of excessive computational requirements or because of lack of a mathematical model. Some of these methods are infinite horizon versions of the ones to be given in Chapter 6 for finite horizon problems, while

others involve ideas that we will cover only briefly in Chapter 6 (see Section 6.3.4).

The presentation in this chapter is original in that it uses the stochastic shortest path problem as the starting point for the analysis of the other problems. This line of development not only explains intuitively the connections between the various types of problems, but also leads to new solution methods. For example, an alternative value iteration algorithm for the average cost problem, based on the connection with the stochastic shortest path problem, is given in the author's paper [Ber98b], and in Section 5.3 of Vol. II. On the other hand, there are important results for undiscounted and average cost problems that cannot be obtained through the connection with the stochastic shortest path problem. Some of these alternative lines of analysis are pursued in Vol. II.

E X E R C I S E S

5.1

A tennis player has a Fast serve and a Slow serve, denoted F and S , respectively. The probability of F (or S) landing in bounds is p_F (or p_S , respectively). The probability of winning the point assuming the serve landed in bounds is q_F (or q_S , respectively). We assume that $p_F < p_S$ and $q_F > q_S$. The problem is to find the serve to be used at each possible scoring situation during a single game in order to maximize the probability of winning that game.

- (a) Formulate this as a stochastic shortest path problem, argue that Assumption 5.2.1 of Section 5.2 holds, and write Bellman's equation.
- (b) Computer assignment: Assume that $q_F = 0.6$, $q_S = 0.4$, and $p_S = 0.95$. Use value iteration to calculate and plot the probability of the server winning a game with optimal serve selection as a function of p_F .

5.2

A quarterback can choose between running and passing the ball on any given play. The number of yards gained by running is integer and is Poisson distributed with parameter λ_r . A pass is incomplete with probability p , is intercepted with probability q , and is completed with probability $1 - p - q$. When completed, a pass gains an integer number of yards that is Poisson distributed with parameter λ_p . We assume that the probability of scoring a touchdown on a single play starting i yards from the goal is equal to the probability of gaining a number of yards greater than or equal to i . We assume also that yardage cannot be lost on any play and that there are no penalties. The ball is turned over to the other team on a fourth down or when an interception occurs.

- (a) Formulate the problem as a stochastic shortest path problem, argue that Assumption 5.2.1 of Section 5.2 holds, and write Bellman's equation.
- (b) Computer assignment: Use value iteration to compute the quarterback's play-selection policy that maximizes the probability of scoring a touchdown on any single drive for $\lambda_r = 3$, $\lambda_p = 10$, $p = 0.4$, and $q = 0.05$.

5.3

A computer manufacturer can be in one of two states. In state 1 his product sells well, while in state 2 his product sells poorly. While in state 1 he can advertise his product in which case the one-stage reward is 4 units, and the transition probabilities are $p_{11} = 0.8$ and $p_{12} = 0.2$. If in state 1, he does not advertise, the reward is 6 units and the transition probabilities are $p_{11} = p_{12} = 0.5$. While in state 2, he can do research to improve his product, in which case the one-stage reward is -5 units, and the transition probabilities are $p_{21} = 0.7$ and $p_{22} = 0.3$. If in state 2 he does not do research, the reward is -3 , and the transition probabilities are $p_{21} = 0.4$ and $p_{22} = 0.6$. Consider the infinite horizon, discounted version of this problem.

- (a) Show that when the discount factor α is sufficiently small, the computer manufacturer should follow the "shortsighted" policy of not advertising (not doing research) while in state 1 (state 2). By contrast, when α is sufficiently close to unity, he should follow the "farsighted" policy of advertising (doing research) while in state 1 (state 2).
- (b) For $\alpha = 0.9$ calculate the optimal policy using policy iteration.
- (c) For $\alpha = 0.99$, use a computer to solve the problem by value iteration, with and without the error bounds (5.25).

5.4

An energetic salesman works every day of the week. He can work in only one of two towns A and B on each day. For each day he works in town A (or B) his expected reward is r_A (or r_B , respectively). The cost for changing towns is c . Assume that $c > r_A > r_B$ and that there is a discount factor $\alpha < 1$.

- (a) Show that for α sufficiently small, the optimal policy is to stay in the town he starts in, and that for α sufficiently close to 1, the optimal policy is to move to town A (if not starting there) and stay in A for all subsequent times.
- (b) Solve the problem for $c = 3$, $r_A = 2$, $r_B = 1$, and $\alpha = 0.9$ using policy iteration.
- (c) Use a computer to solve the problem of part (b) by value iteration, with and without the error bounds (5.25).

5.5

A person has an umbrella that she takes from home to office and vice versa. There is a probability p of rain at the time she leaves home or office independently of earlier weather. If the umbrella is in the place where she is and it rains, she takes the umbrella to go to the other place (this involves no cost). If there is no umbrella and it rains, there is a cost W for getting wet. If the umbrella is in the place where she is but it does not rain, she may take the umbrella to go to the other place (this involves an inconvenience cost V) or she may leave the umbrella behind (this involves no cost). Costs are discounted at a factor $\alpha < 1$.

- (a) Formulate this as an infinite horizon total cost discounted problem. *Hint:* Try to use as few states as possible.
- (b) Characterize the optimal policy as best as you can.

5.6

For the tennis player's problem (Exercise 5.1), show that it is optimal (regardless of score) to use F on both serves if

$$(p_F q_F)/(p_S q_S) > 1,$$

to use S on both serves if

$$(p_F q_F)/(p_S q_S) < 1 + p_F - p_S,$$

and to use F on the first serve and S on the second otherwise.

5.7

Consider the value iteration method for the discounted version of the manufacturer's problem (Example 5.4.2).

$$\begin{aligned} J_{k+1}(i) &= \min \left[K + \alpha(1-p)J_k(0) + \alpha p J_k(1), \right. \\ &\quad \left. ci + \alpha(1-p)J_k(i) + \alpha p J_k(i+1) \right], \quad i = 0, 1, \dots, n-1, \\ J_{k+1}(n) &= K + \alpha(1-p)J_k(0) + \alpha p J_k(1), \end{aligned}$$

where $J_0(i) = 0$ for all i . Show by induction that $J_k(i)$ is monotonically nondecreasing in i .

5.8 www

Consider the policy iteration algorithm for the discounted version of the manufacturer's problem (Example 5.4.2).

- (a) Show that if we start the algorithm with a threshold policy, every subsequently generated policy will be a threshold policy. *Note:* This requires a careful argument.
- (b) Carry out the algorithm for the case $c = 1$, $K = 5$, $n = 10$, $p = 0.5$, $\alpha = 0.9$, and an initial policy that always processes the unfilled orders.

5.9

Solve the average cost version of the computer manufacturer's problem of Exercise 5.3 by using value iteration and by using policy iteration.

5.10

An unemployed worker receives a job offer at each time period, which she may accept or reject. The offered salary takes one of n possible values w^1, \dots, w^n with given probabilities, independently of preceding offers. If she accepts the offer, she must keep the job for the rest of her life at the same salary level. If she rejects the offer, she receives unemployment compensation c for the current period and is eligible to accept future offers. Assume that income is discounted by a factor $\alpha < 1$.

- (a) Show that there is a threshold \bar{w} such that it is optimal to accept an offer if and only if its salary is larger than \bar{w} , and characterize \bar{w} .
- (b) Consider the variant of the problem where there is a given probability p_i that the worker will be fired from her job at any one period if her salary is w^i . Show that the result of part (a) holds in the case where p_i is the same for all i . Analyze the case where p_i depends on i .

5.11

Do part (b) of Exercise 5.10 for the case where income is not discounted and the worker maximizes her average income per period.

5.12 www

Show that one can always take $m = n$ in Assumption 5.2.1. *Hint:* For any π and i , let $S_k(i)$ be the set of states that are reachable with positive probability from i under π in k stages or less. Show that under Assumption 5.2.1, we cannot have $S_k(i) = S_{k+1}(i)$ while $t \neq S_k(i)$.

5.13

Show the error bounds (5.18). These bounds constitute a generalization to the stochastic shortest path problem of the bounds (5.25) for the discounted problem, which have a long history, starting with the work of McQueen [McQ66]; for more recent work see Hansen [Han17], [HaA15]. *Hint:* Complete the details of the following argument. Let $\mu^k(i)$ attain the minimum in the value iteration (5.17) for all i . Then, in vector form, we have

$$J_{k+1} = g_k + P_k J_k,$$

where J_k and g_k are the vectors with components $J_k(i)$, $i = 1, \dots, n$, and $g_k(i, \mu^k(i))$, $i = 1, \dots, n$, respectively, and P_k is the matrix whose components

are the transition probabilities $p_{ij}(\mu^k(i))$. Also from Bellman's equation, we have

$$J^* \leq g_k + P_k J^*,$$

where the vector inequality above is meant to hold separately for each component. Let $e = (1, \dots, 1)'$. Using the above two relations, we have

$$J^* - J_k \leq J^* - J_{k+1} + \bar{c}_k e \leq P_k(J^* - J_k) + \bar{c}_k e. \quad (5.28)$$

Multiplying this relation with P_k and adding $\bar{c}_k e$, we obtain

$$P_k(J^* - J_k) + \bar{c}_k e \leq P_k^2(J^* - J_k) + \bar{c}_k(I + P_k)e.$$

Similarly continuing, we have for all $r \geq 1$

$$J^* - J_{k+1} + \bar{c}_k e \leq P_k^r(J^* - J_k) + \bar{c}_k(I + P_k + \dots + P_k^{r-1})e.$$

For $s = 1, 2, \dots$, the i th component of the vector $P_k^s e$ is equal to the probability $P\{x_s \neq t \mid x_0 = i, \mu^k\}$ that t has not been reached after s stages starting from i and using the stationary policy μ^k . Thus, Assumption 5.2.1 implies that $\lim_{r \rightarrow \infty} P_k^r = 0$, while we have

$$\lim_{r \rightarrow \infty} (I + P_k + \dots + P_k^{r-1})e = N^k,$$

where N^k is the vector $(N^k(1), \dots, N^k(n))'$. Combining the above two relations, we obtain

$$J^* \leq J_{k+1} + \bar{c}_k(N^k - e),$$

proving the desired upper bound.

The lower bound is proved similarly, by using in place of μ^k , an optimal stationary policy μ^* . In particular, in place of Eq. (5.28), we can show that

$$J_k - J^* \leq J_{k+1} - J^* - \underline{c}_k e \leq P^*(J_k - J^*) - \underline{c}_k e,$$

where P^* is the matrix with elements $p_{ij}(\mu^*(i))$. We similarly obtain for all $r \geq 1$

$$J_{k+1} - J^* - \underline{c}_k e \leq (P^*)^r(J_k - J^*) - \underline{c}_k(I + P^* + \dots + (P^*)^{r-1})e,$$

from which $J_{k+1} + \underline{c}_k(N^* - e) \leq J^*$, where N^* is the vector $(N^*(1), \dots, N^*(n))'$.

5.14

5.15

5.16

Consider a problem of operating a machine that can be in any one of n states, denoted $1, 2, \dots, n$. We denote by $g(i)$ the operating cost per period when the machine is in state i , and we assume that

$$g(1) \leq g(2) \leq \dots \leq g(n).$$

The implication here is that state i is better than state $i + 1$, and state 1 corresponds to a machine being in the best condition. The transition probabilities during one period of operation satisfy

$$p_{i(i+1)} > 0 \quad \text{if } i < n,$$

$$p_{ij} = 0 \quad \text{if } j \neq i, j \neq i + 1.$$

We assume that at the start of each period we know the state of the machine and we must choose one of the following two options:

- (1) Let the machine operate one more period in the state it currently is.
- (2) Repair the machine and bring it to the best state 1 at a cost R .

We assume that the machine, once repaired, is guaranteed to stay in state 1 for one period. In subsequent periods, it may deteriorate to states $j > 1$.

- (a) Assume an infinite horizon and a discount factor $\alpha \in (0, 1)$, and show that there is an optimal policy which is a threshold policy; i.e., it takes the form

$$\text{replace if and only if } i \geq i^*,$$

where i^* is some integer.

- (b) Show that the policy iteration method, when started with a threshold policy, generates a sequence of threshold policies.

5.17

5.18

5.19

5.20

5.21

5.22

A treasure hunter has obtained a lease to search a site that contains n treasures, and wants to find a searching policy that maximizes his expected gain over an infinite number of days. At each day, knowing the current number of treasures not yet found, he may decide to continue searching for more treasures at a cost c per day, or to permanently stop searching. If he searches on a day when there are i treasures on the site, he finds $m \in [0, i]$ treasures with given probability $p(m | i)$, where we assume that $p(0 | i) < 1$ for all $i \geq 1$, and that the expected number of treasures found,

$$r(i) = \sum_{m=0}^i mp(m | i),$$

is monotonically nondecreasing with i . Each found treasure is worth 1 unit.

- Formulate the problem as an infinite horizon DP problem.
- Write Bellman's equation. How do you know that this equation holds and has a unique solution?
- Start policy iteration with the policy that never searches. How many policy iterations does it take to find an optimal policy, and what is that optimal policy?

5.23

The latest slot machine model has three arms, labeled 1, 2, and 3. A single play with arm i , where $i = 1, 2, 3$, costs c_i dollars, and has two possible outcomes: a "win," which occurs with probability p_i , and a "loss," which occurs with probability $1 - p_i$. The slot machine pays you m dollars each time you complete a sequence of three successive "wins," with each win obtained using a different arm.

- Consider the problem of finding the arm-playing order that minimizes the expected cost if you are restricted to stop at the first time the machine pays you. Formulate this problem as a stochastic shortest path problem where arm-playing orders are identified with stationary policies, and write Bellman's equation for each stationary policy.
- Show that the expected cost of the arm-playing order ABC is

$$\frac{c_A + p_{ACB} + p_{APBCC} - p_{APBPC}m}{p_{APBPC}}.$$

Show that it is optimal to play the arms in order of decreasing $c_i/(1 - p_i)$.

5.24**5.25**

You have just bought your first car, which raises the issue of where to park it. At the beginning of each day you may either park it in a garage, which costs G per day, or on the street for free. However, in the latter case, you run the risk of getting a parking ticket, which costs T , with probability p_j , where j is the number of consecutive days that the car has been parked on the street (e.g., on the first day you park on the street, you have probability p_1 of getting a ticket, on the second successive day you park on the street, you have probability p_2 , etc). Assume that p_j is monotonically nondecreasing in j , and that you may receive at most one ticket per day when parked on the street. Assume also that there exists an integer m such that $p_m T > G$.

- (a) Formulate this as an infinite horizon discounted cost problem with finite state space and write the corresponding Bellman's equation.
- (b) Characterize as best as you can the optimal policy.
- (c) Let n be the total number of states. Show how to use policy iteration so that it terminates after no more than n iterations. *Hint:* Use threshold policies as in Exercise 5.8.

5.26

An engineer has invented a better mouse trap and is interested in selling it for the right price. At the beginning of each period, he receives a sale offer that takes one of the values s_1, \dots, s_n with corresponding probabilities p_1, \dots, p_n , independently of prior offers. If he accepts the offer he retires from engineering. If he refuses the offer, he may accept subsequent offers but he also runs the risk that a competitor will invent an even better mouse trap, rendering his own unsaleable; this happens with probability $\beta > 0$ at each time period, independently of earlier time periods. While he is overtaken by the competitor, at each time period, he may choose to retire from engineering, or he may choose to invest an amount $v \geq 0$, in which case he has a probability γ to improve his mouse trap, overtake his competitor, and start receiving offers as earlier. The problem is to determine the engineer's strategy to maximize his discounted expected payoff (minus investment cost), assuming a discount factor $\alpha < 1$.

- (a) Formulate the problem as an infinite horizon discounted cost problem and write the corresponding Bellman's equation.
- (b) Characterize as best as you can an optimal policy.
- (c) Assume that there is no discount factor. Does the problem make sense as an average cost per stage problem?
- (d) Assume that there is no discount factor and that the investment cost v is equal to 0. Does the problem make sense as a stochastic shortest path problem, and what is then the optimal policy?

5.27 (Eliminating Self-Transitions)

Consider a stochastic shortest path problem with termination state t , the nontermination states $1, \dots, n$, transition probabilities $p_{ij}(u)$, and expected costs per stage $g(i, u)$. Let Assumption 5.2.1 hold.

- (a) Modify the costs and transition probabilities as follows:

$$\tilde{g}(i, u) = \frac{g(i, u)}{1 - p_{ii}(u)}, \quad i = 1, \dots, n, \quad u \in U(i),$$

$$\tilde{p}_{ij}(u) = \begin{cases} 0 & \text{if } j = i, \\ \frac{p_{ij}(u)}{1 - p_{ii}(u)} & \text{if } j \neq i, \end{cases} \quad i = 1, \dots, n, \quad j = 1, \dots, n, t, \quad u \in U(i),$$

to obtain another stochastic shortest path problem without self-transitions. Show that this modified problem is equivalent to the original in the sense that its stationary policies and optimal policies have the same cost functions. What is the interpretation of the transitions of the modified problem in terms of transitions of the original?

- (b) Fix a policy μ . Let $\{J_k\}$ and $\{\tilde{J}_k\}$ be the sequences of cost vectors generated by value iteration (for the fixed policy) in the original and the modified stochastic shortest path problem, respectively, starting from the same initial vector J_0 . Show that value iteration is faster for the modified problem in the sense that if $J_0 \leq J_1$, then $J_k \leq \tilde{J}_k \leq J^*$ for all k , and if $J_0 \geq J_1$, then $J_k \geq \tilde{J}_k \geq J^*$ for all k .

5.28 (Total Cost Problems with Nonnegative Costs)

This is a theoretical problem whose purpose is to provide some additional analysis for undiscounted cost problems with nonnegative cost per stage, including an extension of the results of Section 5.2 for stochastic shortest path problems. The idea is to use the analysis of Section 5.4 for discounted problems to derive the basic results for total undiscounted cost problems under the assumption that the stage costs are nonnegative and the optimal costs are finite. These results apply, among others, to some stochastic shortest path problems where not all stationary policies are proper and Assumption 5.2.1 is violated. Such problems are discussed in much greater detail in Vol. II; see also the paper by Bertsekas and Yu [BeY16].

Consider a controlled Markov chain with states $i = 1, \dots, n$, controls u chosen from a finite constraint set $U(i)$ for each state i , and transition probabilities $p_{ij}(u)$. (The states may include a cost-free and absorbing termination state, but this is not relevant for the following analysis.) The cost of the k th stage at state i when control u is applied has the form

$$\alpha^k g(i, u), \quad i = 1, \dots, n, \quad u \in U(i),$$

where α is a scalar from $(0, 1]$. Our key assumption is that

$$0 \leq g(i, u), \quad i = 1, \dots, n, \quad u \in U(i).$$

For any policy π , let $J_{\pi,\alpha}$ be the cost function for the α -discounted problem ($\alpha < 1$), and let J_π be the cost function for the problem where $\alpha = 1$. Note that for $\alpha = 1$, we may have $J_\pi(i) = \infty$ for some π and i . However, the limit defining $J_\pi(i)$ exists either as a real number or ∞ , thanks to the assumption $0 \leq g(i, u)$ for all i and u . Let $J_\alpha^*(i)$ and $J^*(i)$ be the optimal costs starting from i , when $\alpha < 1$ and $\alpha = 1$, respectively. We assume that

$$J^*(i) < \infty, \quad i = 1, \dots, n, \quad (5.29)$$

(this is true in particular for the case of a stochastic shortest path problem if there exists a proper stationary policy, i.e., a policy under which there is a positive transition probability path from every state to the termination state).

(a) Show that for all $\alpha < 1$, we have

$$0 \leq J_\alpha^*(i) \leq J^*(i), \quad i = 1, \dots, n.$$

(b) Show that for any admissible policy π , we have

$$\lim_{\alpha \uparrow 1} J_{\pi,\alpha}(i) = J_\pi(i), \quad i = 1, \dots, n.$$

Furthermore,

$$\lim_{\alpha \uparrow 1} J_\alpha^*(i) = J^*(i), \quad i = 1, \dots, n.$$

Hint: To show the first equality, note that for any $\alpha < 1$, N , and $\pi = \{\mu_0, \mu_1, \dots\}$, we have

$$J_\pi(i) \geq J_{\pi,\alpha}(i) \geq \sum_{k=0}^{N-1} \alpha^k E\{g(i_k, \mu_k(i_k)) \mid i_0 = i, \pi\}.$$

Take the limit as $\alpha \rightarrow 1$ and then take the limit as $N \rightarrow \infty$. For the second equality, consider a stationary policy μ and a sequence $\{a_m\} \subset (0, 1)$ with $\alpha_m \rightarrow 1$ such that $J_{\mu,\alpha_m} = J_{\alpha_m}^*$ for all m .

(c) Use Bellman's equation for $\alpha < 1$, to show that J^* satisfies Bellman's equation for $\alpha = 1$:

$$J^*(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J^*(j) \right], \quad i = 1, \dots, n.$$

(d) Let \tilde{J} be such that $0 \leq \tilde{J}(i) < \infty$ for all i . Show that if

$$\tilde{J}(i) \geq \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) \tilde{J}(j) \right], \quad i = 1, \dots, n,$$

then $\tilde{J}(i) \geq J^*(i)$ for all i . Show also that if for some stationary policy μ , we have

$$\tilde{J}(i) \geq g(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) \tilde{J}(j), \quad i = 1, \dots, n,$$

then $\tilde{J}(i) \geq J_\mu(i)$ for all i . *Hint:* Argue that

$$\tilde{J}(i) \geq \min_{u \in U(i)} \left[g(i, u) + \alpha \sum_{j=1}^n p_{ij}(u) \tilde{J}(j) \right], \quad i = 1, \dots, n,$$

use value iteration to show that $\tilde{J} \geq J_\alpha^*$, and take the limit as $\alpha \rightarrow 1$.

(e) For $\alpha = 1$, show that if

$$\mu^*(i) = \arg \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J^*(j) \right], \quad i = 1, \dots, n,$$

then μ^* is optimal. *Hint:* Use part (d) with $\tilde{J} = J^*$.

(f) For $\alpha = 1$, show that for the value iteration method, given by

$$J_{k+1}(i) = \min_{u \in U(i)} \left[g(i, u) + \sum_{j=1}^n p_{ij}(u) J_k(j) \right], \quad i = 1, \dots, n,$$

we have $J_k(i) \rightarrow J^*(i)$, $i = 1, \dots, n$, assuming that

$$0 \leq J_0(i) \leq J^*(i), \quad i = 1, \dots, n.$$

Give examples showing what may happen when this last assumption is violated. *Hint:* Prove the result by first assuming that J_0 is the zero function.

(g) Use the assumption (5.29) to show that the set of states $Z = \{i \mid J^*(i) = 0\}$ is nonempty. Furthermore, under an optimal stationary policy μ^* , the set of states Z is cost-free and absorbing, i.e., $g(i, \mu^*(i)) = 0$ and $p_{ij}(\mu^*(i)) = 0$ for all $i \in Z$ and $j \notin Z$. In addition, μ^* is proper in the sense that for every state $i \notin Z$, under μ^* , there is a positive probability path that starts at i and ends at a state of Z .

5.29 (A Counterexample for Value and Policy Iteration)

Consider a deterministic shortest path problem with a single state/node 1, in addition to the termination state t . At state 1 we can either stay at 1 with cost 0 or move to t with cost b . Note that the set of solutions of Bellman's equation $J(1) = \min \{b, J(1)\}$ is the interval $(-\infty, b]$, and includes $J^*(1)$, which is $J^*(1) = \min\{b, 0\}$. Verify that if $b > 0$, value iteration does not converge to $J^*(1)$, except if started at $J^*(1)$. Moreover, if $b < 0$, policy iteration may oscillate between the optimal and the suboptimal policy. *Note:* This is an example of a semicontractive problem, similar to the linear-quadratic Example 3.1.1. For an analysis of such problems, see the author's abstract DP monograph [Ber13a] and paper [Ber15b].

References

- [ABC65] Atkinson, R. C., Bower, G. H., and Crothers, E. J., 1965. *An Introduction to Mathematical Learning Theory*, Wiley, N. Y.
- [ABF93] Arapostathis, A., Borkar, V., Fernandez-Gaucherand, E., Ghosh, M., and Marcus, S., 1993. “Discrete-Time Controlled Markov Processes with Average Cost Criterion: A Survey,” *SIAM J. on Control and Optimization*, Vol. 31, pp. 282-344.
- [ABG49] Arrow, K. J., Blackwell, D., and Girshick, M. A., 1949. “Bayes and Minimax Solutions of Sequential Design Problems,” *Econometrica*, Vol. 17, pp. 213-244.
- [AGK77] Athans, M., Ku, R., and Gershwin, S. B., 1977. “The Uncertainty Threshold Principle,” *IEEE Trans. on Automatic Control*, Vol. AC-22, pp. 491-495.
- [AHM51] Arrow, K. J., Harris, T., and Marschack, J., 1951. “Optimal Inventory Policy,” *Econometrica*, Vol. 19, pp. 250-272.
- [AKS58] Arrow, K. J., Karlin, S., and Scarf, H., 1958. *Studies in the Mathematical Theory of Inventory and Production*, Stanford Univ. Press, Stanford, CA.
- [Abr90] Abramson, B., 1990. “Expected-Outcome: A General Model of Static Evaluation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 12, pp. 182-193.
- [AdG86] Adams, M., and Guillemin, V., 1986. *Measure Theory and Probability*, Wadsworth and Brooks, Monterey, CA.
- [AnM79] Anderson, B. D. O., and Moore, J. B., 1979. *Optimal Filtering*, Prentice-Hall, Englewood Cliffs, N. J.
- [AoL69] Aoki, M., and Li, M. T., 1969. “Optimal Discrete-Time Control Systems with Cost for Observation,” *IEEE Trans. Automatic Control*, Vol. AC-14, pp. 165-175.
- [AsW94] Åström, K. J., and Wittenmark, B., 1994. *Adaptive Control*, 2nd Edition, Prentice-Hall, Englewood Cliffs, N. J.
- [Ash70] Ash, R. B., 1970. *Basic Probability Theory*, Wiley, N. Y.
- [Ash72] Ash, R. B., 1972. *Real Analysis and Probability*, Academic Press, N. Y.
- [Ast83] Åström, K. J., 1983. “Theory and Applications of Adaptive Control – A Survey,” *Automatica*, Vol. 19, pp. 471-486.
- [AtF66] Athans, M., and Falb, P., 1966. *Optimal Control*, McGraw-Hill, N. Y.
- [BBC11] Bertsimas, D., Brown, D. B., and Caramanis, C., 2011. “Theory and Applications of Robust Optimization,” *SIAM Review*, Vol. 53, pp. 464-501.

- [BBD10] Busoniu, L., Babuska, R., De Schutter, B., and Ernst, D., 2010. Reinforcement Learning and Dynamic Programming Using Function Approximators, CRC Press, N. Y.
- [BBG13] Bertazzi, L., Bosco, A., Guerriero, F., and Lagana, D., 2013. "A Stochastic Inventory Routing Problem with Stock-Out," *Transportation Research, Part C*, Vol. 27, pp. 89-107.
- [BBM17] Borelli, F., Bemporad, A., and Morari, M., 2017. Predictive Control for Linear and Hybrid Systems, Cambridge Univ. Press, Cambridge, UK.
- [BCN16] Bottou, L., Curtis, F. E., and Nocedal, J., 2016. "Optimization Methods for Large-Scale Machine Learning," arXiv preprint arXiv:1606.04838.
- [BGM95] Bertsekas, D. P., Guerriero, F., and Musmanno, R., 1995. "Parallel Shortest Path Methods for Globally Optimal Trajectories," *High Performance Computing: Technology, Methods, and Applications*, (J. Dongarra et al., Eds.), Elsevier.
- [BGM96] Bertsekas, D. P., Guerriero, F., and Musmanno, R., 1996. "Parallel Label Correcting Methods for Shortest Paths," *J. Optimization Theory Appl.*, Vol. 88, 1996, pp. 297-320.
- [BKM05] de Boer, P. T., Kroese, D. P., Mannor, S., and Rubinstein, R. Y. 2005. "A Tutorial on the Cross-Entropy Method," *Annals of Operations Research*, Vol. 134, pp. 19-67.
- [BMO14] Boyd, S., Mueller, M. T., O'Donoghue, B., and Wang, Y., 2014. "Performance Bounds and Suboptimal Policies for Multi-Period Investment," *Foundations and Trends in Optimization*, Vol. 1, pp. 1-72.
- [BMS99] Boltysanski, V., Martini, H., and Soltan, V., 1999. *Geometric Methods and Optimization Problems*, Kluwer, Boston.
- [BNO03] Bertsekas, D. P., Nedić, A., and Ozdaglar, A. E., 2003. *Convex Analysis and Optimization*, Athena Scientific, Belmont, MA.
- [BTW97] Bertsekas, D. P., Tsitsiklis, J. N., and Wu, C., 1997. "Rollout Algorithms for Combinatorial Optimization," *Heuristics*, Vol. 3, pp. 245-262.
- [BYB94] Bradtke, S. J., Ydstie, B. E., and Barto, A. G., 1994. "Adaptive Linear Quadratic Control Using Policy Iteration," *Proc. IEEE American Control Conference*, Vol. 3, pp. 3475-3479.
- [BaB95] Basar, T., and Bernhard, P., 1995. *H_∞ Optimal Control and Related Minimax Design Problems: A Dynamic Game Approach*, Birkhuser, Boston, MA.
- [Bai93] Baird, L. C., 1993. "Advantage Updating," Report WL-TR-93-1146, Wright Patterson AFB, OH.
- [Bai94] Baird, L. C., 1994. "Reinforcement Learning in Continuous Time: Advantage Updating," *International Conf. on Neural Networks*, Orlando, Fla.
- [Bar81] Bar-Shalom, Y., 1981. "Stochastic Dynamic Programming: Caution and Probing," *IEEE Trans. on Automatic Control*, Vol. AC-26, pp. 1184-1195.
- [Bas91] Basar, T., 1991. "Optimum Performance Levels for Minimax Filters, Predictors, and Smoothers," *Systems and Control Letters*, Vol. 16, pp. 309-317.
- [Bas00] Basar, T., 2000. "Risk-Averse Designs: From Exponential Cost to Stochastic Games," In T. E. Djaferis and I. C. Schick, (Eds.), *System Theory: Modeling, Analysis and Control*, Kluwer, Boston, pp. 131-144.
- [BeC99] Bertsekas, D. P., and Castanon, D. A., 1999. "Rollout Algorithms for Stochastic Scheduling Problems," *Heuristics*, Vol. 5, pp. 89-108.

- [BeC04] Bertsekas, D. P., and Castanon, D. A., 2004. Unpublished Collaboration.
- [BeC08] Besse, C., and Chaib-draa, B., 2008. "Parallel Rollout for Online Solution of DEC-POMDPs," Proc. of 21st International FLAIRS Conference, pp. 619-624.
- [BeD62] Bellman, R., and Dreyfus, S., 1962. Applied Dynamic Programming, Princeton Univ. Press, Princeton, N. J.
- [BeG92] Bertsekas, D. P., and Gallager, R. G., 1992. Data Networks, 2nd Edition, Prentice-Hall, Englewood Cliffs, N. J.
- [BeL14] Beyme, S., and Leung, C., 2014. "Rollout Algorithm for Target Search in a Wireless Sensor Network," 80th Vehicular Technology Conference (VTC2014), IEEE, pp. 1-5.
- [BeI96] Bertsekas, D. P., and Ioffe, S., 1996. "Temporal Differences-Based Policy Iteration and Applications in Neuro-Dynamic Programming," Lab. for Info. and Decision Systems Report LIDS-P-2349, Massachusetts Institute of Technology.
- [BeN01] Ben-Tal, A., and Nemirovski, A., 2001. Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications, SIAM, Phila., PA
- [BeP03] Bertsimas, D., and Popescu, I., 2003. "Revenue Management in a Dynamic Network Environment," Transportation Science, Vol. 37, pp. 257-277.
- [BeR71a] Bertsekas, D. P., and Rhodes, I. B., 1971. "On the Minimax Reachability of Target Sets and Target Tubes," Automatica, Vol. 7, pp. 233-247.
- [BeR71b] Bertsekas, D. P., and Rhodes, I. B., 1971. "Recursive State Estimation for a Set-Membership Description of the Uncertainty," IEEE Trans. Automatic Control, Vol. AC-16, pp. 117-128.
- [BeR73] Bertsekas, D. P., and Rhodes, I. B., 1973. "Sufficiently Informative Functions and the Minimax Feedback Control of Uncertain Dynamic Systems," IEEE Trans. Automatic Control, Vol. AC-18, pp. 117-124.
- [BeS78] Bertsekas, D. P., and Shreve, S. E., 1978. Stochastic Optimal Control: The Discrete Time Case, Academic Press, N. Y.; republished by Athena Scientific, Belmont, MA, 1996 (can be downloaded from the author's website).
- [BeT89] Bertsekas, D. P., and Tsitsiklis, J. N., 1989. Parallel and Distributed Computation: Numerical Methods, Prentice-Hall, Englewood Cliffs, N. J.; republished by Athena Scientific, Belmont, MA, 1997.
- [BeT91] Bertsekas, D. P., and Tsitsiklis, J. N., 1991. "An Analysis of Stochastic Shortest Path Problems," Math. Operations Res., Vol. 16, pp. 580-595.
- [BeT96] Bertsekas, D. P., and Tsitsiklis, J. N., 1996. Neuro-Dynamic Programming, Athena Scientific, Belmont, MA.
- [BeT97] Bertsimas, D., and Tsitsiklis, J. N., 1997. Introduction to Linear Optimization, Athena Scientific, Belmont, MA.
- [BeT00] Bertsekas, D. P., and Tsitsiklis, J. N., 2000. "Gradient Convergence of Gradient Methods with Errors," SIAM J. on Optimization, Vol. 36, pp. 627-642.
- [BeT08] Bertsekas, D. P., and Tsitsiklis, J. N., 2008. Introduction to Probability, 2nd Edition, Athena Scientific, Belmont, MA.
- [BeY09] Bertsekas, D. P., and Yu, H., 2009. "Projected Equation Methods for Approximate Solution of Large Linear Systems," J. of Computational and Applied Mathematics, Vol. 227, pp. 27-50.

- [BeY16] Bertsekas, D. P., and Yu, H., 2016. “Stochastic Shortest Path Problems Under Weak Conditions,” Lab. for Information and Decision Systems Report LIDS-2909, MIT.
- [Bel57] Bellman, R., 1957. *Dynamic Programming*, Princeton University Press, Princeton, N. J.
- [Ber70] Bertsekas, D. P., 1970. “On the Separation Theorem for Linear Systems, Quadratic Criteria, and Correlated Noise,” Unpublished Report, Electronic Systems Lab., Massachusetts Institute of Technology.
- [Ber71] Bertsekas, D. P., 1971. “Control of Uncertain Systems With a Set-Membership Description of the Uncertainty,” Ph.D. Dissertation, Massachusetts Institute of Technology, Cambridge, MA (available from the author’s website).
- [Ber72a] Bertsekas, D. P., 1972. “Infinite Time Reachability of State Space Regions by Using Feedback Control,” *IEEE Trans. Automatic Control*, Vol. AC-17, pp. 604-613.
- [Ber72b] Bertsekas, D. P., 1972. “On the Solution of Some Minimax Control Problems,” *Proc. 1972 IEEE Decision and Control Conf.*, New Orleans, LA.
- [Ber74] Bertsekas, D. P., 1974. “Necessary and Sufficient Conditions for Existence of an Optimal Portfolio,” *J. Econ. Theory*, Vol. 8, pp. 235-247.
- [Ber75] Bertsekas, D. P., 1975. “Convergence of Discretization Procedures in Dynamic Programming,” *IEEE Trans. Automatic Control*, Vol. AC-20, pp. 415-419.
- [Ber76] Bertsekas, D. P., 1976. *Dynamic Programming and Stochastic Control*, Academic Press, N. Y.
- [Ber82] Bertsekas, D. P., 1982. “Distributed Dynamic Programming,” *IEEE Trans. Automatic Control*, Vol. AC-27, pp. 610-616.
- [Ber91] Bertsekas, D. P., 1991. *Linear Network Optimization: Algorithms and Codes*, M.I.T. Press, Cambridge, MA.
- [Ber93] Bertsekas, D. P., 1993. “A Simple and Fast Label Correcting Algorithm for Shortest Paths,” *Networks*, Vol. 23, pp. 703-709.
- [Ber97] Bertsekas, D. P., 1997. “Differential Training of Rollout Policies,” *Proc. of the 35th Allerton Conference on Communication, Control, and Computing*, Allerton Park, Ill.
- [Ber98a] Bertsekas, D. P., 1998. *Network Optimization: Continuous and Discrete Models*, Athena Scientific, Belmont, MA.
- [Ber98b] Bertsekas, D. P., 1998. “A New Value Iteration Method for the Average Cost Dynamic Programming Problem,” *SIAM J. on Control and Optimization*, Vol. 36, pp. 742-759.
- [Ber05a] Bertsekas, D. P., 2005. “Dynamic Programming and Suboptimal Control: A Survey from ADP to MPC,” *European J. of Control*, Vol. 11, pp. 310-334.
- [Ber05b] Bertsekas, D. P., 2005. “Rollout Algorithms for Constrained Dynamic Programming,” LIDS Report 2646, MIT.
- [Ber07] Bertsekas, D. P., 2007. “Separable Dynamic Programming and Approximate Decomposition Methods,” *IEEE Trans. on Aut. Control*, Vol. 52, pp. 911-916.
- [Ber09] Bertsekas, D. P., 2009. *Convex Optimization Theory*, Athena Scientific, Belmont, MA.
- [Ber11] Bertsekas, D. P., 2011. “Approximate Policy Iteration: A Survey and Some New Methods,” *J. of Control Theory and Applications*, Vol. 9, pp. 310-335.

- [Ber13a] Bertsekas, D. P., 2013. *Abstract Dynamic Programming*, Athena Scientific, Belmont, MA; a 2nd edition is scheduled to appear in 2017.
- [Ber13b] Bertsekas, D. P., 2013. “Rollout Algorithms for Discrete Optimization: A Survey,” *Handbook of Combinatorial Optimization*, Springer.
- [Ber14] Bertsekas, D. P., 2014. “Robust Shortest Path Planning and Semicontractive Dynamic Programming,” Lab. for Information and Decision Systems Report LIDS-P-2915, MIT; arXiv preprint arXiv:1608.01670; to appear in *Naval Research Logistics*.
- [Ber15a] Bertsekas, D. P., 2015. *Convex Optimization Algorithms*, Athena Scientific, Belmont, MA.
- [Ber15b] Bertsekas, D. P., 2015. “Regular Policies in Abstract Dynamic Programming,” Lab. for Information and Decision Systems Report LIDS-P-3173, MIT; arXiv preprint arXiv:1609.03115.
- [Ber15c] Bertsekas, D. P., 2015. “Value and Policy Iteration in Deterministic Optimal Control and Adaptive Dynamic Programming,” arXiv preprint arXiv:1507.01026; to appear in *IEEE Transactions on Neural Networks and Learning Systems*.
- [Ber16a] Bertsekas, D. P., 2016. *Nonlinear Programming*, 3rd Edition, Athena Scientific, Belmont, MA.
- [Ber16b] Bertsekas, D. P., 2016. “Affine Monotonic and Risk-Sensitive Models in Dynamic Programming,” Lab. for Information and Decision Systems Report LIDS-3204, MIT; arXiv preprint arXiv:1608.01393.
- [BiL97] Birge, J. R., and Louveaux, 1997. *Introduction to Stochastic Programming*, Springer, New York, N. Y.
- [Bis95] Bishop, C. M., 1995. *Neural Networks for Pattern Recognition*, Oxford University Press, N. Y.
- [BIT00] Blondel, V. D., and Tsitsiklis, J. N., 2000. “A Survey of Computational Complexity Results in Systems and Control,” *Automatica*, Vol. 36, pp. 1249-1274.
- [Bla99] Blanchini, F., 1999. “Set Invariance in Control – A Survey,” *Automatica*, Vol. 35, pp. 1747-1768.
- [BoV79] Borkar, V., and Varaiya, P. P., 1979. “Adaptive Control of Markov Chains, I: Finite Parameter Set,” *IEEE Trans. Automatic Control*, Vol. AC-24, pp. 953-958.
- [BoV04] Boyd, S., and Vandenbergue, L., 2004. *Convex Optimization*, Cambridge Univ. Press, Cambridge, U.K.
- [CFH05] Chang, H. S., Hu, J., Fu, M. C., and Marcus, S. I., 2005. “An Adaptive Sampling Algorithm for Solving Markov Decision Processes,” *Operations Research*, Vol. 53, pp. 126-139.
- [CFH13] Chang, H. S., Hu, J., Fu, M. C., and Marcus, S. I., 2013. *Simulation-Based Algorithms for Markov Decision Processes*, (2nd Ed.), Springer, N. Y.
- [CFH16] Chang, H. S., Hu, J., Fu, M. C., and Marcus, S. I., 2016. “Google DeepMind’s AlphaGo,” *ORMS Today, Informs*, Vol. 43.
- [CGC04] Chang, H. S., Givan, R. L., and Chong, E. K. P., 2004. “Parallel Rollout for Online Solution of Partially Observable Markov Decision Processes,” *Discrete Event Dynamic Systems*, Vol. 14, pp. 309-341.
- [CHH02] Campbell, M., Hoane, A. J. and Hsu, F. H., 2002. “Deep Blue,” *Artificial Intelligence*, Vol. 134, pp. 57-83.

- [CaB04] Camacho, E. F., and Bordons, C., 2004. *Model Predictive Control*, 2nd Edition, Springer, New York, N. Y.
- [Cao07] Cao, X. R., 2007. *Stochastic Learning and Optimization: A Sensitivity-Based Approach*, Springer, N. Y.
- [ChT89] Chow, C.-S., and Tsitsiklis, J. N., 1989. "The Complexity of Dynamic Programming," *J. of Complexity*, Vol. 5, pp. 466-488.
- [ChT91] Chow, C.-S., and Tsitsiklis, J. N., 1991. "An Optimal One-Way Multigrid Algorithm for Discrete-Time Stochastic Control," *IEEE Trans. on Automatic Control*, Vol. AC-36, 1991, pp. 898-914.
- [ChV12] Chacon, A., and Vladimirovsky, A., 2012. "Fast Two-Scale Methods for Eikonal Equations," *SIAM J. on Scientific Computing*, Vol. 34, pp. A547-A578.
- [ChV13] Chacon, A., and Vladimirovsky, A., 2013. "A Parallel Heap-Cell Method for Eikonal Equations," *arXiv preprint arXiv:1306.4743*.
- [ChV15] Chacon, A., and Vladimirovsky, A., 2015. "A Parallel Two-Scale Method for Eikonal Equations," *SIAM J. on Scientific Computing*, Vol. 37, pp. A156-A180.
- [Che72] Chernoff, H., 1972. "Sequential Analysis and Optimal Design," *Regional Conference Series in Applied Mathematics*, SIAM, Philadelphia, PA.
- [CoL55] Coddington, E. A., and Levinson, N., 1955. *Theory of Ordinary Differential Equations*, McGraw-Hill, N. Y.
- [Cou06] Coulom, R., 2006. "Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search," *International Conference on Computers and Games*, Springer, pp. 72-83.
- [Cyb89] Cybenko, 1989. "Approximation by Superpositions of a Sigmoidal Function," *Math. of Control, Signals, and Systems*, Vol. 2, pp. 303-314.
- [DeG70] DeGroot, M. H., 1970. *Optimal Statistical Decisions*, McGraw-Hill, N. Y.
- [DeP84] Deo, N., and Pang, C., 1984. "Shortest Path Problems: Taxonomy and Annotation," *Networks*, Vol. 14, pp. 275-323.
- [DeR79] Denardo, E. V., and Rothblum, U. G., 1979. "Optimal Stopping, Exponential Utility, and Linear Programming," *Math. Programming*, Vol. 16, pp. 228-244.
- [Del89] Deller, J. R., 1989. "Set Membership Identification in Digital Signal Processing," *IEEE ASSP Magazine*, Oct., pp. 4-20.
- [DoS80] Doshi, B., and Shreve, S., 1980. "Strong Consistency of a Modified Maximum Likelihood Estimator for Controlled Markov Chains," *J. of Applied Probability*, Vol. 17, pp. 726-734.
- [Dre65] Dreyfus, S. D., 1965. *Dynamic Programming and the Calculus of Variations*, Academic Press, N. Y.
- [Dre69] Dreyfus, S. D., 1969. "An Appraisal of Some Shortest-Path Algorithms," *Operations Research*, Vol. 17, pp. 395-412.
- [Eck68] Eckles, J. E., 1968. "Optimum Maintenance with Incomplete Information," *Operations Research*, Vol. 16, pp. 1058-1067.
- [Elm78] Elmaghraby, S. E., 1978. *Activity Networks: Project Planning and Control by Network Models*, Wiley-Interscience, N. Y.
- [FGL13] Festa, P., Guerriero, F., Lagana, D. and Musmanno, R., 2013. "Solving the Shortest Path Tour Problem," *European J. of Operational Research*, Vol. 230, pp. 464-474.

- [FYG06] Fern, A., Yoon, S. and Givan, R., 2006. "Approximate Policy Iteration with a Policy Language Bias: Solving Relational Markov Decision Processes," *J. of Artificial Intelligence Research*, Vol. 25, pp. 75-118.
- [Fal87] Falcone, M., 1987. "A Numerical Approach to the Infinite Horizon Problem of Deterministic Control Theory," *Appl. Math. Opt.*, Vol. 15, pp. 1-13.
- [FeM94] Fernandez-Gaucherand, E., and Marcus, S. I., 1994. "Risk Sensitive Optimal Control of Hidden Markov Models," *Proc. 33rd IEEE Conf. Dec. Control*, Lake Buena Vista, Fla.
- [FeV02] Ferris, M. C., and Voelker, M. M., 2002. "Neuro-Dynamic Programming for Radiation Treatment Planning," *Numerical Analysis Group Research Report NA-02/06*, Oxford University Computing Laboratory, Oxford University.
- [FeV04] Ferris, M. C., and Voelker, M. M., 2004. "Fractionation in Radiation Treatment Planning," *Mathematical Programming B*, Vol. 102, pp. 387-413.
- [Fel68] Feller, W., 1968. *An Introduction to Probability Theory and its Applications*, Wiley, N. Y.
- [Fei16] Feinberg, E. A., 2016. "Optimality Conditions for Inventory Control," *INFORMS Tutorials in Operations Research*, Online, pp. 14-45.
- [Fis70] Fishburn, P. C., 1970. *Utility Theory for Decision Making*, Wiley, N. Y.
- [For56] Ford, L. R., Jr., 1956. "Network Flow Theory," *Report P-923*, The Rand Corporation, Santa Monica, CA.
- [For73] Forney, G. D., 1973. "The Viterbi Algorithm," *Proc. IEEE*, Vol. 61, pp. 268-278.
- [Fox71] Fox, B. L., 1971. "Finite State Approximations to Denumerable State Dynamic Programs," *J. Math. Anal. Appl.*, Vol. 34, pp. 665-670.
- [Fun89] Funahashi, K., 1989. "On the Approximate Realization of Continuous Mappings by Neural Networks," *Neural Networks*, Vol. 2, pp. 183-192.
- [GBC16] Goodfellow, I., Bengio, J., and Courville, A., *Deep Learning*, MIT Press, Cambridge, MA.
- [GGS13] Gabillon, V., Ghavamzadeh, M., and Scherrer, B., 2013. "Approximate Dynamic Programming Finally Performs Well in the Game of Tetris," in *Advances in Neural Information Processing Systems*, pp. 1754-1762.
- [GTO15] Goodson, J. C., Thomas, B. W., and Ohlmann, J. W., 2015. "Restocking-Based Rollout Policies for the Vehicle Routing Problem with Stochastic Demand and Duration Limits," *Transportation Science*, Vol. 50, pp. 591-607.
- [GaP88] Gallo, G., and Pallottino, S., 1988. "Shortest Path Algorithms," *Annals of Operations Research*, Vol. 7, pp. 3-79.
- [Gal13] Gallager, R. G., 2013. *Stochastic Processes*, Cambridge Univ. Press.
- [GoR85] Gonzalez, R., and Rofman, E., 1985. "On Deterministic Control Problems: An Approximation Procedure for the Optimal Cost, Parts I, II," *SIAM J. Control Optimization*, Vol. 23, pp. 242-285.
- [GoS84] Goodwin, G. C., and Sin, K. S. S., 1984. *Adaptive Filtering, Prediction, and Control*, Prentice-Hall, Englewood Cliffs, N. J.
- [Gos03] Gosavi, A., 2003. *Simulation-Based Optimization: Parametric Optimization Techniques and Reinforcement Learning*, Springer, N. Y.

- [GrA66] Groen, G. J., and Atkinson, R. C., 1966. "Models for Optimizing the Learning Process," *Psychol. Bull.*, Vol. 66, pp. 309-320.
- [Gre11] Grewal, M. S., 2011. *Kalman Filtering*, Springer, Berlin.
- [GuF63] Gunckel, T. L., and Franklin, G. R., 1963. "A General Solution for Linear Sampled-Data Control," *Trans. ASME Ser. D. J. Basic Engrg.*, Vol. 85, pp. 197-201.
- [GuM01] Guerriero, F., and Musmanno, R., 2001. "Label Correcting Methods to Solve Multicriteria Shortest Path Problems," *J. Optimization Theory Appl.*, Vol. 111, pp. 589-613.
- [GuM03] Guerriero, F., and Mancini, M., 2003. "A Cooperative Parallel Rollout Algorithm for the Sequential Ordering Problem," *Parallel Computing*, Vol. 29, pp. 663-677.
- [HJG16] Huang, Q., Jia, Q. S., and Guan, X., 2016. "Robust Scheduling of EV Charging Load with Uncertain Wind Power Integration," *IEEE Transactions on Smart Grid*.
- [HMS55] Holt, C. C., Modigliani, F., and Simon, H. A., 1955. "A Linear Decision Rule for Production and Employment Scheduling," *Management Sci.*, Vol. 2, pp. 1-30.
- [HNR68] Hart, P. E., Nilsson, N. J. and Raphael, B., 1968. "A Formal Basis for the Heuristic Determination of Minimum Cost Paths," *IEEE Trans. on Systems Science and Cybernetics*, Vol. 4, pp. 100-107.
- [HNR72] Hart, P. E., Nilsson, N. J. and Raphael, B., 1972. "Correction to a Formal Basis for the Heuristic Determination of Minimum Cost Paths," *ACM SIGART Bulletin*, Vol. 37, pp. 28-29.
- [HPC96] Helmsen, J., Puckett, E. G., Colella, P., and Dorr, M., 1996. "Two New Methods for Simulating Photolithography Development," *SPIE's 1996 International Symposium on Microlithography*, pp. 253-261.
- [HSW89] Hornick, K., Stinchcombe, M., and White, H., 1989. "Multilayer Feedforward Networks are Universal Approximators," *Neural Networks*, Vol. 2, pp. 359-159.
- [HaA15] Hansen, E. A. and Abdoulahi, I., 2015. "Efficient Bounds in Heuristic Search Algorithms for Stochastic Shortest Path Problems," in *AAAI*, pp. 3283-3290.
- [HaL82] Hajek, B., and van Loon, T., 1982. "Decentralized Dynamic Control of a Multiaccess Broadcast Channel," *IEEE Trans. Automatic Control*, Vol. AC-27, pp. 559-569.
- [Hak70] Hakansson, N. H., 1970. "Optimal Investment and Consumption Strategies under Risk for a Class of Utility Functions," *Econometrica*, Vol. 38, pp. 587-607.
- [Hak71] Hakansson, N. H., 1971. "On Myopic Portfolio Policies, With and Without Serial Correlation of Yields," *The J. of Business of the University of Chicago*, Vol. 44, pp. 324-334.
- [Han80] Hansen, P., 1980. "Bicriterion Path Problems," in *Multiple-Criteria Decision Making: Theory and Applications*, Edited by G. Fandel and T. Gal, Springer Verlag, Heidelberg, Germany, pp. 109-127.
- [Han06] Hansen, N., 2006. "The CMA Evolution Strategy: A Comparing Review," in *Towards a New Evolutionary Computation*, Springer Berlin Heidelberg, pp. 75-102.
- [Han17] Hansen, E. A., 2017. "Error Bounds for Stochastic Shortest Path Problems," *Math. of Operations Research*, forthcoming.
- [Hay09] Haykin, S., 2009. *Neural Networks and Learning Machines*, (3rd Edition), Prentice-Hall, Englewood-Cliffs, N. J.

- [Hes66] Hestenes, M. R., 1966. *Calculus of Variations and Optimal Control Theory*, Wiley, N. Y.
- [Her89] Hernandez-Lerma, O., 1989. *Adaptive Markov Control Processes*, Springer, N. Y.
- [HoK71] Hoffman, K., and Kunze, R., 1971. *Linear Algebra*, 2nd ed., Prentice-Hall, Englewood Cliffs, N. J.
- [IEE71] IEEE Trans. Automatic Control, 1971. Special Issue on Linear-Quadratic Gaussian Problem, Vol. AC-16.
- [IoS96] Ioannou, P. A., and Sun, J., 1996. *Robust Adaptive Control*, Prentice-Hall, Englewood Cliffs, N. J.
- [JBE94] James, M. R., Baras, J. S., and Elliott, R. J., 1994. "Risk-Sensitive Control and Dynamic Games for Partially Observed Discrete-Time Nonlinear Systems," IEEE Trans. on Automatic Control, Vol. AC-39, pp. 780-792.
- [Jac73] Jacobson, D. H., 1973. "Optimal Stochastic Linear Systems With Exponential Performance Criteria and their Relation to Deterministic Differential Games," IEEE Trans. Automatic Control, Vol. AC-18, pp. 124-131.
- [Jaf84] Jaffe, J. M., 1984. "Algorithms for Finding Paths with Multiple Constraints," Networks, Vol. 14, pp. 95-116.
- [Jew63] Jewell, W., 1963. "Markov Renewal Programming I and II," Operations Research, Vol. 2, pp. 938-971.
- [JoT61] Joseph, P. D., and Tou, J. T., 1961. "On Linear Control Theory," AIEE Trans., Vol. 80 (II), pp. 193-196.
- [Jon90] Jones, L. K., 1990. "Constructive Approximations for Neural Networks by Sigmoidal Functions," Proceedings of the IEEE, Vol. 78, pp. 1586-1589.
- [KGB82] Kimemia, J., Gershwin, S. B., and Bertsekas, D. P., 1982. "Computation of Production Control Policies by a Dynamic Programming Technique," in *Analysis and Optimization of Systems*, A. Bensoussan and J. L. Lions (eds.), Springer, N. Y., pp. 243-269.
- [KKK95] Krstic, M., Kanellakopoulos, I., Kokotovic, P., 1995. *Nonlinear and Adaptive Control Design*, J. Wiley, N. Y.
- [KLB92] Kosut, R. L., Lau, M. K., and Boyd, S. P., 1992. "Set-Membership Identification of Systems with Parametric and Nonparametric Uncertainty," IEEE Trans. on Automatic Control, Vol. AC-37, pp. 929-941.
- [KLC98] Kaelbling, L.P., Littman, M. L., and Cassandra, A. R., 1998. "Planning and Acting in Partially Observable Stochastic Domains," Artificial Intelligence, Vol. 101, pp. 99-134.
- [KaD66] Karush, W., and Dear, E. E., 1966. "Optimal Stimulus Presentation Strategy for a Stimulus Sampling Model of Learning," J. Math. Psychology, Vol. 3, pp. 15-47.
- [KaK58] Kalman, R. E., and Koepcke, R. W., 1958. "Optimal Synthesis of Linear Sampling Control Systems Using Generalized Performance Indexes," Trans. ASME, Vol. 80, pp. 1820-1826.
- [KaW94] Kall, P., and Wallace, S. W., 1994. *Stochastic Programming*, Wiley, Chichester, UK.
- [Kal60] Kalman, R. E., 1960. "A New Approach to Linear Filtering and Prediction Problems," Trans. ASME Ser. D. J. Basic Engrg., Vol. 82, pp. 35-45.

- [KeG88] Keerthi, S. S., and Gilbert, E. G., 1988. "Optimal, Infinite Horizon Feedback Laws for a General Class of Constrained Discrete Time Systems: Stability and Moving-Horizon Approximations," *J. Optimization Theory Appl.*, Vol. 57, pp. 265-293.
- [Kim82] Kimemia, J., 1982. "Hierarchical Control of Production in Flexible Manufacturing Systems," Ph.D. Thesis, Dep. of Electrical Engineering and Computer Science, Massachusetts Institute of Technology.
- [KoC15] Kouvaritakis, B., and Cannon, M., 2015. *Model Predictive Control: Classical, Robust and Stochastic*, Springer, N. Y.
- [KoS06] Kocsis, L., and Szepesvari, C., 2006. "Bandit Based Monte-Carlo Planning," *Proc. of 17th European Conference on Machine Learning*, Berlin, pp. 282-293.
- [KuA77] Ku, R., and Athans, M., 1977. "Further Results on the Uncertainty Threshold Principle," *IEEE Trans. on Automatic Control*, Vol. AC-22, pp. 866-868.
- [KuD92] Kushner, H. J., and Dupuis, P. G., 1992. *Numerical Methods for Stochastic Control Problems in Continuous Time*, Springer, N. Y.
- [KuL82] Kumar, P. R., and Lin, W., 1982. "Optimal Adaptive Controllers for Unknown Markov Chains," *IEEE Trans. Automatic Control*, Vol. AC-27, pp. 765-774.
- [KuV86] Kumar, P. R., and Varaiya, P. P., 1986. *Stochastic Systems: Estimation, Identification, and Adaptive Control*, Prentice-Hall, Englewood Cliffs, N. J.
- [KuV97] Kurzhanski, A., and Valyi, I., 1997. *Ellipsoidal Calculus for Estimation and Control*, Birkhauser, Boston, MA.
- [Kum83] Kumar, P. R., 1983. "Optimal Adaptive Control of Linear-Quadratic-Gaussian Systems," *SIAM J. on Control and Optimization*, Vol. 21, pp. 163-178.
- [Kum85] Kumar, P. R., 1985. "A Survey of Some Results in Stochastic Adaptive Control," *SIAM J. on Control and Optimization*, Vol. 23, pp. 329-380.
- [LGW16] Lan, Y., Guan, X., and Wu, J., 2016. "Rollout Strategies for Real-Time Multi-Energy Scheduling in Microgrid with Storage System," *IET Generation, Transmission and Distribution*, Vol. 10, pp. 688-696.
- [LLL08] Lewis, F. L., Liu, D., and Lendaris, G. G., 2008. Special Issue on Adaptive Dynamic Programming and Reinforcement Learning in Feedback Control, *IEEE Trans. on Systems, Man, and Cybernetics, Part B*, Vol. 38, No. 4.
- [LLP93] Leshno, M., Lin, V. Y., Pinkus, A., and Schocken, S., 1993. "Multilayer Feed-forward Networks with a Nonpolynomial Activation Function can Approximate any Function," *Neural Networks*, Vol. 6, pp. 861-867.
- [LWW17] Liu, D., Wei, Q., Wang, D., and Yang, X., 2017. *Adaptive Dynamic Programming with Applications in Optimal Control*, Springer, Berlin.
- [Las85] Lasserre, J. B., 1985. "A Mixed Forward-Backward Dynamic Programming Method Using Parallel Computation," *J. Optimization Theory Appl.*, Vol. 45, pp. 165-168.
- [LeL13] Lewis, F. L., and Liu, D., (Eds), 2013. *Reinforcement Learning and Approximate Dynamic Programming for Feedback Control*, Wiley, Hoboken, N. J.
- [Lev84] Levy, D., 1984. *The Chess Computer Handbook*, B. T. Batsford Ltd., London.
- [LiR71] Lippman, S. A., and Ross, S. M., 1971. "The Streetwalker's Dilemma: A Job-Shop Model," *SIAM J. of Appl. Math.*, Vol. 20, pp. 336-342.

- [LiR06] Lincoln, B., and Rantzer, A., 2006. "Relaxing Dynamic Programming," *IEEE Trans. Automatic Control*, Vol. 51, pp. 1249-1260.
- [LiW15] Li, H. and Womer, N. K., 2015. "Solving Stochastic Resource-Constrained Project Scheduling Problems by Closed-Loop Approximate Dynamic Programming," *European J. of Operational Research*, Vol. 246, pp. 20-33.
- [LjS83] Ljung, L., and Soderstrom, T., 1983. *Theory and Practice of Recursive Identification*, MIT Press, Cambridge, MA.
- [Lov91a] Lovejoy, W. S., 1991. "Computationally Feasible Bounds for Partially Observed Markov Decision Processes," *Operations Research*, Vol. 39, pp. 1621-175.
- [Lov91b] Lovejoy, W. S., 1991. "A Survey of Algorithmic Methods for Partially Observed Markov Decision Processes," *Annals of Operations Research*, Vol. 18, pp. 47-66.
- [Lue69] Luenberger, D. G., 1969. *Optimization by Vector Space Methods*, Wiley, N. Y.
- [Lue84] Luenberger, D. G., 1984. *Linear and Nonlinear Programming*, Addison-Wesley, Reading, MA.
- [MHK98] Meuleau, N., Hauskrecht, M., Kim, K.-E., Peshkin, L., Kaelbling, L. K., and Dean, T., 1998. "Solving Very Large Weakly Coupled Markov Decision Processes," *Proc. of the Fifteenth National Conference on Artificial Intelligence*, Madison, WI, pp. 165-172.
- [MMB02] McGovern, A., Moss, E., and Barto, A., 2002. "Building a Basic Building Block Scheduler Using Reinforcement Learning and Rollouts," *Machine Learning*, Vol. 49, pp. 141-160.
- [MPP04] Meloni, C., Pacciarelli, D., and Pranzo, M., 2004. "A Rollout Metaheuristic for Job Shop Scheduling Problems," *Annals of Operations Research*, Vol. 131, pp. 215-235.
- [MVZ95] Mulvey, J. M., Vanderbei, R. J., and Zenios, S. A., 1995. "Robust Optimization of Large-Scale Systems," *Operations Research*, Vol. 43, pp. 264-281.
- [MaJ15] Mastin, A., and Jaillet, P., 2015. "Average-Case Performance of Rollout Algorithms for Knapsack Problems," *J. of Optimization Theory and Applications*, Vol. 165, pp. 964-984.
- [Mac02] Maciejowski, J. M., 2002. *Predictive Control with Constraints*, Addison-Wesley, Reading, MA.
- [Mar84] Martins, E. Q. V., 1984. "On a Multicriteria Shortest Path Problem," *European J. of Operational Research*, Vol. 16, pp. 236-245.
- [McQ66] MacQueen, J., 1966. "A Modified Dynamic Programming Method for Markovian Decision Problems," *J. Math. Anal. Appl.*, Vol. 14, pp. 38-43.
- [Mik79] Mikhailov, V. A., 1979. *Methods of Random Multiple Access*, Candidate Engineering Thesis, Moscow Institute of Physics and Technology, Moscow.
- [MoL99] Morari, M., and Lee, J. H., 1999. "Model Predictive Control: Past, Present, and Future," *Computers and Chemical Engineering*, Vol. 23, pp. 667-682.
- [Mos68] Mossin, J., 1968. "Optimal Multi-Period Portfolio Policies," *J. Business*, Vol. 41, pp. 215-229.
- [MuS08] Munos, R., and Szepesvari, C., 2008. "Finite-Time Bounds for Fitted Value Iteration," *J. of Machine Learning Research*, Vol. 1, pp. 815-857.
- [Mun14] Munos, R., 2014. "From Bandits to Monte-Carlo Tree Search: The Optimistic Principle Applied to Optimization and Planning," *Foundations and Trends in Machine Learning*, Vol. 7, pp. 1-129.

- [NeW88] Nemhauser, G. L., and Wolsey, L. A., 1988. *Integer and Combinatorial Optimization*, Wiley, N. Y.
- [New75] Newborn, M., 1975. *Computer Chess*, Academic Press, N. Y.
- [Nic66] Nicholson, T., 1966. "Finding the Shortest Route Between Two Points in a Network," *The Computer Journal*, Vol. 9, pp. 275-280.
- [Nil71] Nilsson, N. J., 1971. *Problem-Solving Methods in Artificial Intelligence*, McGraw-Hill, N. Y.
- [Nil80] Nilsson, N. J., 1980. *Principles of Artificial Intelligence*, Morgan-Kaufmann, San Mateo, Ca.
- [OBP16] Osband, I., Blundell, C., Pritzel, A., and Van Roy, B., 2016. "Deep Exploration Via Bootstrapped DQN," *Advances in Neural Information Processing Systems*, Vol. 29.
- [OVW16] Osband, I., Van Roy, B., and Wen, Z., 2016. "Generalization and Exploration Via Randomized Value Functions," *Proc. of the 33rd International Conference on Machine Learning*, pp. 2377-2386.
- [Osb16] Osband, I., 2016. *Deep Exploration via Randomized Value Functions*, Ph.D. Dissertation, Stanford University.
- [PBG65] Pontryagin, L. S., Boltyanski, V., Gamkrelidze, R., and Mishchenko, E., 1965. *The Mathematical Theory of Optimal Processes*, Interscience Publishers, Inc., N. Y.
- [PBT98] Polymenakos, L. C., Bertsekas, D. P., and Tsitsiklis, J. N., 1998. "Efficient Algorithms for Continuous-Space Shortest Path Problems," *IEEE Trans. on Automatic Control*, Vol. 43, pp. 278-283.
- [PaT87] Papadimitriou, C. H., and Tsitsiklis, J. N., 1987. "The Complexity of Markov Decision Processes," *Math. Operations Res.*, Vol. 12, pp. 441-450.
- [Pap74] Pape, V., 1974. "Implementation and Efficiency of Moore Algorithms for the Shortest Path Problem," *Math. Progr.*, Vol. 7, pp. 212-222.
- [Pat01] Patek, S. D., 2001. "On Terminating Markov Decision Processes with a Risk Averse Objective Function," *Automatica*, Vol. 37, pp. 1379-1386.
- [Pea84] Pearl, J., 1984. *Heuristics*, Addison-Wesley, Reading, MA.
- [Pic90] Picone, J., 1990. "Continuous Speech Recognition Using Hidden Markov Models," *IEEE ASSP Magazine*, July Issue, pp. 26-41.
- [Pin95] Pinedo, M., 1995. *Scheduling: Theory, Algorithms, and Systems*, Prentice-Hall, Englewood Cliffs, N. J.
- [Pos14] Post, H. N., 2014. *The Shortest Path Problem on Real Road Networks: Theory, Algorithms and Computations*, TU Delft, Delft University of Technology.
- [Pow07] Powell, W. B., 2007. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*, J. Wiley and Sons, Hoboken, N. J; a 2nd edition appeared in 2011.
- [PrS94] Proakis, J. G., and Salehi, M., 1994. *Communication Systems Engineering*, Prentice-Hall, Englewood Cliffs, N. J.
- [Pra64] Pratt, J. W., 1964. "Risk Aversion in the Small and in the Large," *Econometrica*, Vol. 32, pp. 300-307.
- [Pre95] Prekopa, A., 1995. *Stochastic Programming*, Kluwer, Boston.
- [RSS12] Runarsson, T. P., Schoenauer, M., and Sebag, M., 2012. "Pilot, Rollout and Monte Carlo Tree Search Methods for Job Shop Scheduling," in *Learning and Intelligent*

Optimization, Springer, Berlin, pp. 160-174.

[Rab89] Rabiner, L. R., 1989. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proc. of the IEEE*, Vol. 77, pp. 257-286.

[Roc70] Rockafellar, R. T., 1970. *Convex Analysis*, Princeton University Press, Princeton, N. J.

[Ros70] Ross, S. M., 1970. *Applied Probability Models with Optimization Applications*, Holden-Day, San Francisco, CA.

[Ros83] Ross, S. M., 1983. *Introduction to Stochastic Dynamic Programming*, Academic Press, N. Y.

[Ros85] Ross, S. M., 1985. *Probability Models*, Academic Press, Orlando, Fla.

[Ros12] Ross, S. M., 2012. *Simulation*, 5th Edition, Academic Press, Orlando, Fla.

[RuK04] Rubinstein, R. Y., and Kroese, D. P., 2004. *The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization*, Springer, N. Y.

[RuK17] Rubinstein, R. Y., and Kroese, D. P., 2017. *Simulation and the Monte Carlo Method*, 3rd Edition, J. Wiley, N. Y.

[Rud76] Rudin, W., 1976. *Real Analysis*, 3rd Edition, McGraw-Hill, N. Y.

[Rus97] Rust, J., 1997. "Using Randomization to Break the Curse of Dimensionality," *Econometrica*, Vol. 65, pp. 487-516.

[SBB89] Sastry, S., Bodson, M., and Bartram, J. F., 1989. *Adaptive Control: Stability, Convergence, and Robustness*, Prentice-Hall, Englewood Cliffs, N. J.

[SBP04] Si, J., Barto, A., Powell, W., and Wunsch, D., (Eds.) 2004. *Learning and Approximate Dynamic Programming*, IEEE Press, N. Y.

[SGC02] Savagaonkar, U., Givan, R., and Chong, E. K. P., 2002. "Sampling Techniques for Zero-Sum, Discounted Markov Games," in *Proc. 40th Allerton Conference on Communication, Control and Computing*, Monticello, Ill.

[SGG15] Scherrer, B., Ghavamzadeh, M., Gabillon, V., Lesner, B., and Geist, M., 2015. "Approximate Modified Policy Iteration and its Application to the Game of Tetris," *J. of Machine Learning Research*, Vol. 16, pp. 1629-1676.

[SHB15] Simroth, A., Holfeld, D., and Brunsch, R., 2015. "Job Shop Production Planning under Uncertainty: A Monte Carlo Rollout Approach," *Proc. of the International Scientific and Practical Conference*, Vol. 3, pp. 175-179.

[SHM16] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., and Dieleman, S., 2016. "Mastering the Game of Go with Deep Neural Networks and Tree Search," *Nature*, Vol. 529, pp. 484-489.

[SLJ13] Sun, B., Luh, P. B., Jia, Q. S., Jiang, Z., Wang, F., and Song, C., 2013. "Building Energy Management: Integrated Control of Active and Passive Heating, Cooling, Lighting, Shading, and Ventilation Systems," *IEEE Trans. on Automation Science and Engineering*, Vol. 10, pp. 588-602.

[SZL08] Sun, T., Zhao, Q., Lun, P., and Tomastik, R., 2008. "Optimization of Joint Replacement Policies for Multipart Systems by a Rollout Framework," *IEEE Transactions on Automation Science and Engineering*, Vol. 5, pp. 609-619.

[Sam69] Samuelson, P. A., 1969. "Lifetime Portfolio Selection by Dynamic Stochastic Programming," *Review of Economics and Statistics*, Vol. 51, pp. 239-246.

- [Sar87] Sargent, T. J., 1987. *Dynamic Macroeconomic Theory*, Harvard Univ. Press, Cambridge, MA.
- [Sca60] Scarf, H., 1960. "The Optimality of (s, S) Policies for the Dynamic Inventory Problem," *Proceedings of the 1st Stanford Symposium on Mathematical Methods in the Social Sciences*, Stanford University Press, Stanford, CA.
- [Sch68] Schweppe, F. C., 1968. "Recursive State Estimation; Unknown but Bounded Errors and System Inputs," *IEEE Trans. Automatic Control*, Vol. AC-13.
- [Sch97] Schaeffer, J., 1997. *One Jump Ahead*, Springer, N. Y.
- [Sch13] Scherrer, B., 2013. "Performance Bounds for Lambda Policy Iteration and Application to the Game of Tetris," *J. of Machine Learning Research*, Vol. 14, pp. 1181-1227.
- [Sec00] Secomandi, N., 2000. "Comparing Neuro-Dynamic Programming Algorithms for the Vehicle Routing Problem with Stochastic Demands," *Computers and Operations Research*, Vol. 27, pp. 1201-1225.
- [Sec01] Secomandi, N., 2001. "A Rollout Policy for the Vehicle Routing Problem with Stochastic Demands," *Operations Research*, Vol. 49, pp. 796-802.
- [Sec03] Secomandi, N., 2003. "Analysis of a Rollout Approach to Sequencing Problems with Stochastic Routing Applications," *J. of Heuristics*, Vol. 9, pp. 321-352.
- [Set99a] Sethian, J. A., 1999. *Level Set Methods and Fast Marching Methods Evolving Interfaces in Computational Geometry, Fluid Mechanics, Computer Vision, and Materials Science*, Cambridge University Press, N. Y.
- [Set99b] Sethian, J. A., 1999. "Fast Marching Methods," *SIAM Review*, Vol. 41, pp. 199-235.
- [Sha50] Shannon, C., 1950. "Programming a Digital Computer for Playing Chess," *Phil. Mag.*, Vol. 41, pp. 356-375.
- [Shi64] Shiryaev, A. N., 1964. "On Markov Sufficient Statistics in Non-Additive Bayes Problems of Sequential Analysis," *Theory of Probability and Applications*, Vol. 9, pp. 604-618.
- [Shi66] Shiryaev, A. N., 1966. "On the Theory of Decision Functions and Control by an Observation Process with Incomplete Data," *Selected Translations in Math. Statistics and Probability*, Vol. 6, pp. 162-188.
- [Shr81] Shreve, S. E., 1981. "A Note on Optimal Switching Between Two Activities," *Naval Research Logistics Quarterly*, Vol. 28, pp. 185-190.
- [Sim56] Simon, H. A., 1956. "Dynamic Programming Under Uncertainty with a Quadratic Criterion Function," *Econometrica*, Vol. 24, pp. 74-81.
- [Sim06] Simon, D., 2006. *Optimal State Estimation: Kalman, H-Infinity, and Nonlinear Approaches*, J. Wiley, N. Y.
- [Skl88] Sklar, B., 1988. *Digital Communications: Fundamentals and Applications*, Prentice-Hall, Englewood Cliffs, N. J.
- [SIL91] Slotine, J.-J. E., and Li, W., *Applied Nonlinear Control*, Prentice-Hall, Englewood Cliffs, N. J.
- [SmS73] Smallwood, R. D., and Sondik, E. J., 1973. "The Optimal Control of Partially Observable Markov Processes Over a Finite Horizon," *Operations Res.*, Vol. 11, pp. 1071-1088.

- [Sma71] Smallwood, R. D., 1971. "The Analysis of Economic Teaching Strategies for a Simple Learning Model," *J. Math. Psychology*, Vol. 8, pp. 285-301.
- [Sob75] Sobel, M. J., 1975. "Ordinal Dynamic Programming," *Management Science*, Vol. 21, pp. 967-975.
- [Son71] Sondik, E. J., 1971. "The Optimal Control of Partially Observable Markov Processes," Ph.D. Dissertation, Department of Engineering-Economic Systems, Stanford University, Stanford, CA.
- [StL89] Stokey, N. L., and Lucas, R. E., 1989. *Recursive Methods in Economic Dynamics*, Harvard University Press, Cambridge, MA.
- [StW91] Stewart, B. S., and White, C. C., 1991. "Multiobjective A^* ," *J. ACM*, Vol. 38, pp. 775-814.
- [Sti94] Stirzaker, D., 1994. *Elementary Probability*, Cambridge University Press, Cambridge.
- [Str65] Striebel, C. T., 1965. "Sufficient Statistics in the Optimal Control of Stochastic Systems," *J. Math. Anal. Appl.*, Vol. 12, pp. 576-592.
- [Str76] Strang, G., 1976. *Linear Algebra and its Applications*, Academic Press, N. Y.
- [SuB98] Sutton, R., and Barto, A. G., 1998. *Reinforcement Learning*, MIT Press, Cambridge, MA.
- [SzL06] Szita, I., and Lorinz, A., 2006. "Learning Tetris Using the Noisy Cross-Entropy Method," *Neural Computation*, Vol. 18, pp. 2936-2941.
- [Sze10] Szepesvari, C., 2010. *Algorithms for Reinforcement Learning*, Morgan and Claypool Publishers, San Francisco, CA.
- [TGL13] Tesauro, G., Gondek, D. C., Lenchner, J., Fan, J., and Prager, J. M., 2013. "Analysis of Watson's Strategies for Playing Jeopardy!," *J. of Artificial Intelligence Research*.
- [TeG96] Tesauro, G., and Galperin, G. R., 1996. "On-Line Policy Improvement Using Monte Carlo Search," presented at the 1996 Neural Information Processing Systems Conference, Denver, CO; also in M. Mozer et al. (eds.), *Advances in Neural Information Processing Systems 9*, MIT Press (1997).
- [ThS09] Thiery, C., and Scherrer, B., 2009. "Improvements on Learning Tetris with Cross Entropy," *International Computer Games Association Journal*, Vol. 32, pp. 23-33.
- [Tes89] Tesauro, G., "Connectionist Learning of Expert Preferences by Comparison Training," in D. Touretzky (Ed.), *Advances in Neural Information Processing Systems*, (NIPS-88), Morgan Kaufmann, San Mateo, CA, pp. 99-106.
- [Tes01] Tesauro, G., "Comparison Training of Chess Evaluation Functions," in J. Furnkranz, M. Kumbat (Eds.), *Machines that Learn to Play Games*, Nova Science Publishers, pp. 117-130.
- [ThS09] Thiery, C., and Scherrer, B., 2009. "Improvements on Learning Tetris with Cross-Entropy," *International Computer Games Association J.*, Vol. 32, pp. 23-33.
- [The54] Theil, H., 1954. "Econometric Models and Welfare Maximization," *Weltwirtsch. Arch.*, Vol. 72, pp. 60-83.
- [TsV96] Tsitsiklis, J. N., and Van Roy, B., 1996. "Feature-Based Methods for Large-Scale Dynamic Programming," *Machine Learning*, Vol. 22, pp. 59-94.

- [Tsi84a] Tsitsiklis, J. N., 1984. "Convexity and Characterization of Optimal Policies in a Dynamic Routing Problem," *J. Optimization Theory Appl.*, Vol. 44, pp. 105-136.
- [Tsi84b] Tsitsiklis, J. N., 1984. "Periodic Review Inventory Systems with Continuous Demand and Discrete Order Sizes," *Management Sci.*, Vol. 30, pp. 1250-1254.
- [Tsi87] Tsitsiklis, J. N., 1987. "Analysis of a Multiaccess Control Scheme," *IEEE Trans. Automatic Control*, Vol. AC-32, pp. 1017-1020.
- [Tsi95] Tsitsiklis, J. N., 1995. "Efficient Algorithms for Globally Optimal Trajectories," *IEEE Trans. Automatic Control*, Vol. AC-40, pp. 1528-1538.
- [TuP03] Tu, F., and Pattipati, K. R., 2003. "Rollout Strategies for Sequential Fault Diagnosis," *IEEE Trans. on Systems, Man and Cybernetics, Part A*, pp. 86-99.
- [VVL13] Vrabie, D., Vamvoudakis, K. G., and Lewis, F. L., 2013. *Optimal Adaptive Control and Differential Games by Reinforcement Learning Principles*, The Institution of Engineering and Technology, London.
- [VaW89] Varaiya, P., and Wets, R. J-B., 1989. "Stochastic Dynamic Optimization Approaches and Computation," *Mathematical Programming: State of the Art*, M. Iri and K. Tanabe (eds.), Kluwer, Boston, pp. 309-332.
- [Vei65] Veinott, A. F., Jr., 1965. "The Optimal Inventory Policy for Batch Ordering," *Operations Res.*, Vol. 13, pp. 424-432.
- [Vei66] Veinott, A. F., Jr., 1966. "The Status of Mathematical Inventory Theory," *Management Sci.*, Vol. 12, pp. 745-777.
- [Vin74] Vincke, P., 1974. "Problemes Multicriteres," *Cahiers du Centre d' Etudes de Recherche Operationnelle*, Vol. 16, pp. 425-439.
- [Vit67] Viterbi, A. J., 1967. "Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm," *IEEE Trans. on Info. Theory*, Vol. IT-13, pp. 260-269.
- [Vla08] Vladimirovsky, A., 2008. "Label-Setting Methods for Multimode Stochastic Shortest Path Problems on Graphs," *Math. of Operations Research*, Vol. 33, pp. 821-838.
- [WCG03] Wu, G., Chong, E. K. P., and Givan, R. L., 2003. "Congestion Control Using Policy Rollout," *Proc. 2nd IEEE CDC*, Maui, Hawaii, pp. 4825-4830.
- [Wal47] Wald, A., 1947. *Sequential Analysis*, Wiley, N. Y.
- [WeB99] Weaver, L., and Baxter, J., 1999. "Reinforcement Learning From State and Temporal Differences," *Tech. Report*, Department of Computer Science, Australian National University.
- [WeP80] Weiss, G., and Pinedo, M., 1980. "Scheduling Tasks with Exponential Service Times on Nonidentical Processors to Minimize Various Cost Functions," *J. Appl. Prob.*, Vol. 17, pp. 187-202.
- [Wen14] Wen, Z., 2014. *Efficient Reinforcement Learning with Value Function Generalization*, Ph.D. Dissertation, Stanford University.
- [WhH80] White, C. C., and Harrington, D. P., 1980. "Application of Jensen's Inequality to Adaptive Suboptimal Design," *J. Optimization Theory Appl.*, Vol. 32, pp. 89-99.
- [WhS89] White, C. C., and Scherer, W. T., 1989. "Solution Procedures for Partially Observed Markov Decision Processes," *Operations Res.*, Vol. 30, pp. 791-797.
- [Whi69] White, D. J., 1969. *Dynamic Programming*, Holden-Day, San Francisco, CA.

- [Whi78] Whitt, W., 1978. "Approximations of Dynamic Programs I," *Math. Operations Res.*, Vol. 3, pp. 231-243.
- [Whi79] Whitt, W., 1979. "Approximations of Dynamic Programs II," *Math. Operations Res.*, Vol. 4, pp. 179-185.
- [Whi82] Whittle, P., 1982. *Optimization Over Time*, Wiley, N. Y., Vol. 1, 1982, Vol. 2, 1983.
- [Whi88] Whittle, P., 1988. "Restless Bandits: Activity Allocation in a Changing World," *J. of Applied Probability*, pp. 287-298.
- [Whi90] Whittle, P., 1990. *Risk-Sensitive Optimal Control*, Wiley, N. Y.
- [Wil71] Willems, J., 1971. "Least Squares Stationary Optimal Control and the Algebraic Riccati Equation," *IEEE Trans. on Automatic Control*, Vol. 16, pp. 621-634.
- [Wit66a] Witsenhausen, H. S., 1966. "Minimax Control of Uncertain Systems," Ph.D. Dissertation, Massachusetts Institute of Technology, Cambridge, MA.
- [Wit66b] Witsenhausen, H. S., 1966. "A Comparison of Closed-Loop and Open-Loop Optimum Systems," *IEEE Trans. Automatic Control*, Vol. AC-11, pp. 619-621.
- [Wit68] Witsenhausen, H. S., 1968. "Sets of Possible States of Linear Systems Given Perturbed Observations," *IEEE Trans. Automatic Control*, Vol. AC-13, pp. 556-558.
- [Wit69] Witsenhausen, H. S., 1969. "Inequalities for the Performance of Suboptimal Uncertain Systems," *Automatica*, Vol. 5, pp. 507-512.
- [Wit70] Witsenhausen, H. S., 1970. "On Performance Bounds for Uncertain Systems," *SIAM J. on Control*, Vol. 8, pp. 55-89.
- [Wit71] Witsenhausen, H. S., 1971. "Separation of Estimation and Control for Discrete-Time Systems," *Proc. IEEE*, Vol. 59, pp. 1557-1566.
- [Wol98] Wolsey, L. A., 1998. *Integer Programming*, Wiley, N. Y.
- [WuB99] Wu, C. C., and Bertsekas, D. P., 1999. "Distributed Power Control Algorithms for Wireless Networks," unpublished report, available from the author's [www](#) site.
- [YDR04] Yan, X., Diaconis, P., Rusmevichientong, P., and Van Roy, B., 2004. "Solitaire: Man Versus Machine," *Advances in Neural Information Processing Systems*, Vol. 17, pp. 1553-1560.
- [YuB04] Yu, H., and Bertsekas, D. P., 2004. "Discretized Approximations for POMDP with Average Cost," *Proc. of 20th Conference on Uncertainty in Artificial Intelligence*, Banff, Canada.
- [YuB12] Yu, H., and Bertsekas, D. P., 2012. "Weighted Bellman Equations and their Applications in Dynamic Programming," *Lab. for Information and Decision Systems Report LIDS-P-2876*, MIT.
- [YuB15] Yu, H., and Bertsekas, D. P., 2015. "A Mixed Value and Policy Iteration Method for Stochastic Control with Universally Measurable Policies," *Math. of Operations Research*, Vol. 40, pp. 926-968.