

*Rollout, Policy Iteration,
and
Distributed Reinforcement Learning*

by

Dimitri P. Bertsekas

Arizona State University
and
Massachusetts Institute of Technology

WWW site for book information and orders

<http://www.athenasc.com>



Athena Scientific, Belmont, Massachusetts

Athena Scientific
Post Office Box 805
Nashua, NH 03060
U.S.A.

Email: info@athenasc.com
WWW: <http://www.athenasc.com>

Cover photography by Dimitri Bertsekas.
Stars over the Stata Center at MIT (built on the location of the old Building
20 where Claude Shannon had his first office as a professor in 1956).

© 2020 Dimitri P. Bertsekas
All rights reserved. No part of this book may be reproduced in any form
by any electronic or mechanical means (including photocopying, recording,
or information storage and retrieval) without permission in writing from
the publisher.

Publisher's Cataloging-in-Publication Data

Bertsekas, Dimitri P.
Rollout, Policy Iteration, and Distributed Reinforcement Learning
Includes Bibliography and Index
1. Mathematical Optimization. 2. Dynamic Programming. I. Title.
QA402.5 .B465 2020 519.703 00-91281

ISBN-10: 1-886529-07-8, ISBN-13: 978-1-886529-07-6

ABOUT THE AUTHOR

Dimitri Bertsekas studied Mechanical and Electrical Engineering at the National Technical University of Athens, Greece, and obtained his Ph.D. in system science from the Massachusetts Institute of Technology. He has held faculty positions with the Engineering-Economic Systems Department, Stanford University, and the Electrical Engineering Department of the University of Illinois, Urbana. Since 1979 he has been teaching at the Electrical Engineering and Computer Science Department of the Massachusetts Institute of Technology (M.I.T.), where he is McAfee Professor of Engineering. In 2019, he joined the School of Computing, Informatics, and Decision Systems Engineering at the Arizona State University, Tempe, AZ, as Fulton Professor of Computational Decision Making.

Professor Bertsekas' teaching and research have spanned several fields, including deterministic optimization, dynamic programming and stochastic control, large-scale and distributed computation, artificial intelligence, and data communication networks. He has authored or coauthored numerous research papers and eighteen books, several of which are currently used as textbooks in MIT classes, including "Dynamic Programming and Optimal Control," "Data Networks," "Introduction to Probability," and "Nonlinear Programming."

Professor Bertsekas was awarded the INFORMS 1997 Prize for Research Excellence in the Interface Between Operations Research and Computer Science for his book "Neuro-Dynamic Programming" (co-authored with John Tsitsiklis), the 2001 AACC John R. Ragazzini Education Award, the 2009 INFORMS Expository Writing Award, the 2014 AACC Richard Bellman Heritage Award, the 2014 INFORMS Khachiyan Prize for Lifetime Accomplishments in Optimization, and the 2015 MOS/SIAM George B. Dantzig Prize. In 2018 he shared with his coauthor, John Tsitsiklis, the 2018 INFORMS John von Neumann Theory Prize for the contributions of the research monographs "Parallel and Distributed Computation" and "Neuro-Dynamic Programming." Professor Bertsekas was elected in 2001 to the United States National Academy of Engineering for "pioneering contributions to fundamental research, practice and education of optimization/control theory, and especially its application to data communication networks."

ATHENA SCIENTIFIC
OPTIMIZATION AND COMPUTATION SERIES

1. Rollout, Policy Iteration, and Distributed Reinforcement Learning, by Dimitri P. Bertsekas, 2020, ISBN 978-1-886529-07-6, 376 pages
2. Reinforcement Learning and Optimal Control, by Dimitri P. Bertsekas, 2019, ISBN 978-1-886529-39-7, 388 pages
3. Abstract Dynamic Programming, 2nd Edition, by Dimitri P. Bertsekas, 2018, ISBN 978-1-886529-46-5, 360 pages
4. Dynamic Programming and Optimal Control, Two-Volume Set, by Dimitri P. Bertsekas, 2017, ISBN 1-886529-08-6, 1270 pages
5. Nonlinear Programming, 3rd Edition, by Dimitri P. Bertsekas, 2016, ISBN 1-886529-05-1, 880 pages
6. Convex Optimization Algorithms, by Dimitri P. Bertsekas, 2015, ISBN 978-1-886529-28-1, 576 pages
7. Convex Optimization Theory, by Dimitri P. Bertsekas, 2009, ISBN 978-1-886529-31-1, 256 pages
8. Introduction to Probability, 2nd Edition, by Dimitri P. Bertsekas and John N. Tsitsiklis, 2008, ISBN 978-1-886529-23-6, 544 pages
9. Convex Analysis and Optimization, by Dimitri P. Bertsekas, Angelia Nedić, and Asuman E. Ozdaglar, 2003, ISBN 1-886529-45-0, 560 pages
10. Network Optimization: Continuous and Discrete Models, by Dimitri P. Bertsekas, 1998, ISBN 1-886529-02-7, 608 pages
11. Network Flows and Monotropic Optimization, by R. Tyrrell Rockafellar, 1998, ISBN 1-886529-06-X, 634 pages
12. Introduction to Linear Optimization, by Dimitris Bertsimas and John N. Tsitsiklis, 1997, ISBN 1-886529-19-1, 608 pages
13. Parallel and Distributed Computation: Numerical Methods, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1997, ISBN 1-886529-01-9, 718 pages
14. Neuro-Dynamic Programming, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1996, ISBN 1-886529-10-8, 512 pages
15. Constrained Optimization and Lagrange Multiplier Methods, by Dimitri P. Bertsekas, 1996, ISBN 1-886529-04-3, 410 pages
16. Stochastic Optimal Control: The Discrete-Time Case, by Dimitri P. Bertsekas and Steven E. Shreve, 1996, ISBN 1-886529-03-5, 330 pages

Contents

1. Dynamic Programming Principles	
1.1. Deterministic Dynamic Programming	p. 2
1.1.1. Basic Finite Horizon Problem Formulation	p. 2
1.1.2. The Dynamic Programming Algorithm	p. 5
1.1.3. Approximation in Value Space	p. 7
1.2. Stochastic Dynamic Programming	p. 10
1.2.1. Finite Horizon Problems	p. 10
1.2.2. Infinite Horizon Problems - An Overview	p. 14
1.3. Examples, Variations, and Simplifications	p. 20
1.3.1. Discrete Deterministic Optimization	p. 21
1.3.2. Problems with a Termination State	p. 25
1.3.3. State Augmentation, Time Delays, and Forecasts	p. 29
1.3.4. Partial State Information and Belief States	p. 32
1.4. Reinforcement Learning and Optimal Control - Some Terminology	p. 35
1.5. Notes and Sources	p. 37
2. Rollout and Policy Improvement	
2.1. Approximation in Value and Policy Space	p. 43
2.1.1. Approximation in Value Space - One-Step and Multistep Lookahead	p. 43
2.1.2. Approximation in Policy Space	p. 47
2.1.3. Combined Approximation in Value and Policy Space	p. 49
2.2. General Issues of Approximation in Value Space	p. 53
2.2.1. Model-Based and Model-Free Implementations	p. 53
2.2.2. Off-Line and On-Line Implementations	p. 54
2.2.3. Methods for Cost-to-Go Approximation	p. 56
2.2.4. Methods for Simplification of the Lookahead Minimization	p. 58
2.2.5. Simplification of the Lookahead Minimization by Q-Factor Approximation	p. 59

2.3. Rollout and the Policy Improvement Principle	p. 62
2.3.1. On-Line Rollout for Deterministic Discrete	
Optimization	p. 64
2.3.2. Using Multiple Base Heuristics - Parallel Rollout	p. 72
2.3.3. The Fortified Rollout Algorithm	p. 73
2.3.4. Truncated Rollout with Multistep Lookahead and	
Terminal Cost Approximation	p. 76
2.3.5. Rollout with Small Stage Costs and Long Horizon -	
Continuous-Time Rollout	p. 78
2.3.6. Rollout with an Expert	p. 88
2.4. Stochastic Rollout and Monte Carlo Tree Search	p. 91
2.4.1. Simulation-Based Implementation of the Rollout	
Algorithm	p. 94
2.4.2. Rollout and Monte Carlo Tree Search	p. 98
2.4.3. Randomized Policy Improvement by Monte Carlo	
Tree Search	p. 102
2.4.4. Rollout Parallelization	p. 103
2.4.5. The Effect of Errors in Rollout - Variance	
Reduction	p. 104
2.5. Rollout for Infinite-Spaces Problems - Optimization	
Heuristics	p. 107
2.5.1. Rollout for Infinite-Spaces Deterministic Problems	p. 107
2.5.2. Rollout Based on Stochastic Programming	p. 111
2.6. Notes and Sources	p. 114
3. Specialized Rollout Algorithms	
3.1. Model Predictive Control	p. 120
3.1.1. Target Tubes and Constrained Controllability	
Condition	p. 127
3.1.2. Model Predictive Control with Terminal Cost	p. 131
3.1.3. Variants of Model Predictive Control	p. 132
3.2. Multiagent Rollout	p. 135
3.2.1. Multiagent Coupling Through Constraints	p. 145
3.2.2. Multiagent Rollout for Separable and Multiarmed	
Bandit Problems	p. 147
3.2.3. Multiagent Model Predictive Control	p. 150
3.2.4. Asynchronous Distributed Multiagent Rollout	p. 151
3.3. Constrained Rollout for Deterministic Optimization	p. 155
3.3.1. State-Constrained Rollout and Target Tubes	p. 156
3.3.2. Rollout with Trajectory Constraints	p. 160
3.3.3. Constrained Multiagent Rollout	p. 169
3.4. Constrained Rollout - Combinatorial and Discrete	
Optimization	p. 172
3.4.1. A General Discrete Optimization Problem	p. 172

- 3.4.2. Multidimensional Assignment p. 179
- 3.5. Surrogate Dynamic Programming and Rollout p. 187
 - 3.5.1. Rollout for Bayesian Optimization p. 189
- 3.6. Rollout for Minimax Control p. 193
- 3.7. Notes and Sources p. 196

4. Learning Values and Policies

- 4.1. Approximation Architectures p. 204
 - 4.1.1. Feature-Based Architectures p. 205
 - 4.1.2. Training of Linear and Nonlinear Architectures p. 215
- 4.2. Neural Networks p. 219
 - 4.2.1. Training of Neural Networks p. 223
 - 4.2.2. Multilayer and Deep Neural Networks p. 224
- 4.3. Training of Cost Functions in Approximate DP p. 226
 - 4.3.1. Fitted Value Iteration p. 226
 - 4.3.2. Q-Factor Parametric Approximation p. 228
 - 4.3.3. Advantage Updating - Approximating Q-Factor Differences p. 230
 - 4.3.4. Differential Training of Cost Differences for Rollout p. 233
- 4.4. Training of Policies in Approximate DP p. 235
 - 4.4.1. Perpetual Rollout with Value and Policy Networks - Multiprocessor Parallelization p. 239
- 4.5. Notes and Sources p. 240

5. Infinite Horizon: Distributed and Multiagent Algorithms

- 5.1. Stochastic Shortest Path and Discounted Problems p. 248
- 5.2. Exact and Approximate Policy Iteration p. 260
 - 5.2.1. Policy Iteration and Rollout p. 261
 - 5.2.2. Optimistic and Multistep Policy Iteration - Truncated Rollout p. 265
 - 5.2.3. Policy Iteration for Q-Factors p. 268
 - 5.2.4. Multiagent Rollout and Policy Iteration p. 270
 - 5.2.5. Approximation in Value Space p. 277
 - 5.2.6. Performance Bounds for Truncated Rollout and Approximate Policy Iteration p. 279
- 5.3. Abstract View of Infinite Horizon Problems p. 290
- 5.4. Multiagent Value and Policy Iteration p. 301
 - 5.4.1. Convergence to an Agent-by-Agent Optimal Policy p. 305
 - 5.4.2. Optimistic Multiagent PI Algorithms p. 310
- 5.5. Asynchronous Distributed Value Iteration p. 313
 - 5.5.1. State Space Partitioning p. 314
 - 5.5.2. Asynchronous Convergence Theorem p. 315
- 5.6. Asynchronous Distributed Policy Iteration p. 318
 - 5.6.1. Randomized Asynchronous Optimistic Policy

Iteration p. 320

5.6.2. Asynchronous Optimistic Policy Iteration with a
Uniform Fixed Point p. 323

5.6.3. Approximate Policy Iteration - Asynchronous
Multiprocessor Parallelization p. 330

5.7. Notes and Sources p. 332

References p. 337

Index p. 357

Preface

We know the past but cannot control it. We control the future but cannot know it.

Claude Shannon

In this research monograph we discuss the solution of large and challenging multistage decision problems using methods of reinforcement learning (RL for short), also referred to by other names such as approximate dynamic programming and neuro-dynamic programming. We will focus on a subset of methods which are based on the idea of *policy iteration*, i.e., starting from some policy and generating one or more improved policies.

If just one improved policy is generated, this is called *rollout*, which, based on broad and consistent computational experience, appears to be one of the simplest and most reliable of all RL methods. Rollout is also well-suited for on-line model-free implementation and on-line replanning. Approximate policy iteration can be viewed as repeated application of rollout. This is one of the most prominent types of RL methods. It can be implemented using data generated by the system itself, a process known as *self-learning*, and value and policy approximation architectures, including neural networks.

Approximate policy iteration is more ambitious than rollout, but it is a strictly off-line method, and it is generally far more computationally intensive (of course rollout may also require a lot of on-line computation). This motivates the use of parallel and distributed computation. One of the purposes of the monograph is to discuss distributed (possibly asynchronous) methods that relate to rollout and policy iteration, both in the context of an exact and an approximate implementation involving neural networks or other approximation architectures.

One of the contributions of the monograph is to develop variants of rollout and policy iteration for problems with a multiagent structure, where the control consists of multiple components, each associated with a separate agent. In particular, we introduce a new approach to lookahead

simplification through the use of *multiagent rollout*, which allows the dramatic reduction of the computational requirements for one-step lookahead when the control consists of multiple components, and connects with the theme of distributed asynchronous implementation.

Multiagent rollout also has a strong connection with a well-developed body of research with a long history: the theory of teams and the notion of person-by-person optimality. In particular, we develop an infinite horizon dynamic programming methodology, which includes value and policy iteration methods that converge to a person-by-person optimal policy. While our multiagent schemes are based on fully shared agent information, they are also well suited as a starting point for approximations, in the context of on-line autonomous decision making by multiple agents each coordinating in varying degrees with the other agents. In this context, agent information that is not shared by other agents, is appropriately estimated, with the estimates being treated as if they were exact.

Several of the ideas that we develop in some depth in this monograph have been central in the implementation of recent high profile successes, such as the AlphaZero program for playing chess, Go, and other games. In addition to the fundamental process of successive policy iteration/improvement, this program includes the use of deep neural networks for representation of both value functions and policies, the extensive use of large scale parallelization, and the simplification of lookahead minimization, through methods involving Monte Carlo tree search and pruning of the lookahead tree. In this monograph, we also focus on policy iteration, value and policy neural network representations, parallel and distributed computation, and lookahead simplification. Thus while there are significant differences, the principal design ideas that form the core of this monograph are shared by the AlphaZero architecture, except that we develop these ideas in a broader and less application-specific framework.

Another subject that we deal with in some detail is model predictive control (MPC for short), one of the most prominent control system design methods at present. One of the reasons is that classical forms of MPC are closely related to (and indeed can be viewed as) rollout algorithms, thereby providing a connection with reinforcement learning, which is beneficial in two ways. On one hand the MPC context provides rich crossfertilization opportunities with the analytical and algorithmic ideas of rollout and RL; for example the notion of sequential improvement in rollout is intimately connected to Lyapunov stability analysis in MPC, and the target tube ideas that are central in MPC may prove useful in the context of constrained rollout and policy iteration. On the other hand the dynamic programming and RL methodologies point the way to extensions of MPC based on self-learning, approximate policy iteration, simulation, the treatment of stochastic and set membership uncertainty, and the use of distributed computation.

In our development of several of the topics of this book we rely on

methodology that is covered in greater depth in the 1996 neuro-dynamic programming book [BeT96] (jointly written with J. Tsitsiklis) as well as the author’s recent RL book [Ber19a]. However, we aim to develop rollout and approximate policy iteration beyond the books [BeT96] and [Ber19a]. In particular, we present new research, relating to distributed asynchronous computation, partitioned architectures, and multiagent systems. We also indicate how our methods are well-suited for several types of challenging large scale optimization problems, such as combinatorial/discrete optimization, as well as partially observed Markov decision problems (POMDP).

This monograph took shape in the fall of 2019 and was based on two separate but related lines of research for distributed large-scale computation:

- (a) My work on rollout, policy iteration, and value iteration for multi-agent DP problems, including those arising in discrete deterministic optimization settings [Ber19c], [Ber19d], [Ber20].
- (b) My work on distributed policy iteration algorithms with state space-partitioned architectures. These ideas were extended, implemented, and applied to some large-scale POMDP problems in collaboration with my Arizona State University (ASU) colleagues Sushmita Bhattacharya, Sahil Badyal, Thomas Wheeler, and Stephanie Gil [BBW20]. This work is also connected with my joint research on asynchronous distributed state space-partitioned policy iteration with Huizhen Yu [BeY10], [BeY12], [YuB13], which is presented in Section 5.6 of this monograph.

Most of the book was written while teaching a research-oriented course at ASU starting in January 2020. The hospitable and stimulating environment at ASU contributed much to my productivity during this period, and for this I am very thankful to several colleagues and students, including Stephanie Gil, Giulia Pedrielli, and Petr Sulc, and my teaching assistant, Sushmita Bhattacharya. I have also appreciated fruitful interactions with colleagues and students outside ASU, particularly Yuchao Li, who also provided valuable proofreading support.

Finally, I would like to dedicate this monograph to the creative genius of Claude Shannon, the father of information theory, but also the father of computer chess. His approximate dynamic programming ideas, which predated the work of Bellman, live on inside the AlphaZero program, the most impressive success story of reinforcement learning up to now.

Dimitri P. Bertsekas

July 2020