

Studies in Computational Intelligence

Volume 503

Series Editor

J. Kacprzyk, Warsaw, Poland

For further volumes:

<http://www.springer.com/series/7092>

Todd Hester

TEXPLORE: Temporal Difference Reinforcement Learning for Robots and Time-Constrained Domains

Todd Hester
Department of Computer Science
University of Texas at Austin
Austin, Texas
USA

ISSN 1860-949X ISSN 1860-9503 (electronic)
ISBN 978-3-319-01167-7 ISBN 978-3-319-01168-4 (eBook)
DOI 10.1007/978-3-319-01168-4
Springer Cham Heidelberg New York Dordrecht London

Library of Congress Control Number: 2013942265

© Springer International Publishing Switzerland 2013

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

To my family, especially my lovely wife Beth

Foreword

It is a great pleasure and honor to be able to write the foreword for this book, representing the culmination of Todd Hester’s Ph.D. thesis research at The University of Texas at Austin. It was my good fortune to be Todd’s advisor during his years as a Computer Science graduate student. I therefore was able to participate in and enjoy the adventure of starting with the kernel of an idea and fully developing it into a full-blown dissertation.

Todd’s research is in the area of Reinforcement learning (RL), a machine learning paradigm that focuses on enabling computers and robots to learn to perform sequential tasks. Though grounded in some very theoretically elegant results, and showing great promise for enabling learning-based robots that could be deployed in the real world, most RL algorithms to date require either too much data (training experience) or too much computation to be practical on real-world problems.

This book introduces `TEXPLORE`, one of the first RL algorithms to be both data-efficient and computation-efficient enough to work on real robots in the real world. The core idea of `TEXPLORE` is the realization that scaling up to large domains in real time requires actively reasoning about which states *not* to explore. Most RL algorithms to date still insist on exhaustive exploration: visiting every state, or in continuous settings, every region of the state space. Doing so is necessary if the goal is finding the *optimal* policy because in principle, any unvisited state could be a “gold mine.” However when accepting that perfection can sometimes be the enemy of the good, it becomes clear that exploration must be more focused.

The distinguishing characteristic of `TEXPLORE` is that it learns *in real-time, while continuing to act*. As this book documents fully, it has been demonstrated to learn a speed-control task on a real autonomous vehicle during the course of two minutes of continual driving. In addition to this driving task, the algorithms introduced in this book have been validated on real humanoid robots, and in carefully controlled, simulation environments.

This book also includes an investigation into the application of `TEXPLORE` to the idea of intrinsically motivated RL. Analogous to curiosity on the part of human learners, intrinsically motivated RL requires guiding exploration by properties of the environment, rather than based on external reward.

This book is important for the field for these novel algorithms themselves, and also for the fact that it opens up several exciting directions for future research. By releasing the associated source code as an RL package within the Robot

Operating System (ROS) development environment, Todd has made it easy for future researchers to build upon his contributions.

Overall, for both newcomers to the field, and for practitioners looking for nuanced detail, this book has plenty to offer. Whichever your perspective, I trust you will enjoy reading it!

Austin, Texas
June, 2013

Peter Stone

Preface

This book presents the main results of the research I conducted for my Ph.D. thesis at the University of Texas at Austin. The main focus of the research is on developing new reinforcement learning methods that enable fast and robust learning on robots in real-time.

Robots have the potential to solve many problems in society, because of their ability to work in dangerous places doing necessary jobs that no one wants or is able to do. One barrier to their widespread deployment is that they are mainly limited to tasks where it is possible to hand-program behaviors for every situation that may be encountered. For robots to meet their potential, they need methods that enable them to learn and adapt to novel situations that they were not programmed for. Reinforcement learning (RL) is a paradigm for learning sequential decision making processes and could solve the problems of learning and adaptation on robots. This book identifies four key challenges that must be addressed for an RL algorithm to be practical for robotic control tasks. These *RL for Robotics Challenges* are: 1) it must learn in very few samples; 2) it must learn in domains with continuous state features; 3) it must handle sensor and/or actuator delays; and 4) it should continually select actions in real time. This book focuses on addressing all four of these challenges. In particular, this book is focused on *time-constrained domains* where the first challenge is critically important. In these domains, the agent’s lifetime is not long enough for it to explore the domain thoroughly, and it must learn in very few samples.

Although existing RL algorithms successfully address one or more of the *RL for Robotics Challenges*, no prior algorithm addresses all four of them. To fill this gap, this book introduces **TEXPLORE**, the first algorithm to address all four challenges. **TEXPLORE** is a model-based RL method that learns a random forest model of the domain which generalizes dynamics to unseen states. Each tree in the random forest model represents a hypothesis of the domain’s true dynamics, and the agent uses these hypotheses to explore states that are promising for the final policy, while ignoring states that do not appear promising. With sample-based planning and a novel parallel architecture, **TEXPLORE** can select actions continually in real time whenever necessary.

We empirically evaluate each component of **TEXPLORE** in comparison with other state-of-the-art approaches. In addition, we present modifications of **TEXPLORE**’s exploration mechanism for different types of domains. The key result of this book is a demonstration of **TEXPLORE** learning to control the velocity of an autonomous vehicle on-line, in real time, while running on-board the robot. After

controlling the vehicle for only two minutes, *TEXPLORE* is able to learn to move the pedals of the vehicle to drive at the desired velocities. The work presented in this book represents an important step towards applying RL to robotics and enabling robots to perform more tasks in society. By enabling robots to learn in few actions while acting on-line in real time on robots with continuous state and actuator delays, *TEXPLORE* significantly broadens the applicability of RL to robots.

This book would not have been possible without help from a great number of people. First and foremost, I want to thank my advisor Peter Stone, whose guidance, advice, and support has been invaluable. There are also many other graduate students who helped and collaborated with me on this research. In particular, Nick Jong let me assist on an AAMAS paper on reinforcement learning in my first year as a graduate student and gave me a great start on RL.

Austin, Texas
April, 2013

Todd Hester

Table of Contents

1	Introduction	1
1.1	Time-Constrained Domains	3
1.2	Algorithm Overview	5
1.3	Contributions	7
1.4	Overview	7
2	Background and Problem Specification	11
2.1	Markov Decision Problems	11
2.2	Value Function Reinforcement Learning	13
2.2.1	Model-Free Methods	13
2.2.2	Model-Based Methods	14
2.2.3	Factored Models	17
2.2.4	Planning	18
2.3	Time-Constrained Domains	19
2.4	A Specific Problem	21
2.5	Chapter Summary	23
3	Real Time Architecture	25
3.1	Monte Carlo Tree Search (MCTS) Planning	26
3.1.1	Domains with Delay	29
3.2	Parallel Architecture	31
3.3	Chapter Summary	34
4	The TEXPLORE Algorithm	35
4.1	Model Learning	35
4.1.1	Models of Continuous Domains	38
4.1.2	Domains with Delays	38
4.1.3	Dependent Feature Transitions	41
4.2	Exploration	43
4.3	The Complete TEXPLORE Algorithm	48
4.4	Chapter Summary	49
5	Empirical Evaluation	51
5.1	Challenge 1: Sample Efficiency and Exploration	52
5.1.1	Simulated Vehicle Velocity Control	53
5.1.2	Fuel World	56

5.2	Challenge 2: Modeling Continuous Domains	62
5.2.1	Simulated Vehicle Velocity Control	62
5.2.2	Continuous Task Performance	63
5.3	Challenge 3: Delayed Actions	67
5.3.1	Simulated Vehicle Velocity Control	67
5.3.2	Delayed Gridworld	69
5.4	Challenge 4: Real Time Action	73
5.4.1	Simulated Vehicle Velocity Control	73
5.4.2	Mountain Car	75
5.5	Dependent Transitions	78
5.6	TEXPLORE on a Physical Robot	83
5.7	Chapter Summary	84
6	Further Examination of Exploration	85
6.1	Explicit Exploration	87
6.1.1	Methodology	87
6.1.2	Empirical Evaluation	88
6.2	Variance and Novelty Intrinsic Rewards	94
6.2.1	Methodology	95
6.2.2	Empirical Evaluation	98
6.3	On-Line Learning of Exploration Parameters	105
6.3.1	Methodology	105
6.3.2	Empirical Evaluation	108
6.4	Empirical Comparison	115
6.5	Chapter Summary	119
7	Related Work	121
7.1	Sample Efficiency	121
7.1.1	Exploration	122
7.1.2	Intrinsic Motivation	124
7.1.3	Bayesian Methods	126
7.1.4	Models	128
7.2	Continuous Domains	129
7.3	Observation and Action Delays	130
7.4	Real-Time Architectures	131
7.5	Real-World Problem Domains	133
7.6	Chapter Summary	134
8	Discussion and Conclusion	137
8.1	Summary	137
8.2	Contributions	139
8.3	Discussion	140
8.4	Future Work	142
8.4.1	Expanded Applicability of RL	142
8.4.2	Exploration	144
8.4.3	Opponent Modeling	145

8.4.4 Lifelong Learning	146
8.4.5 Summary	147
8.5 Conclusion	147
A TEXPLORE Pseudo-code	149
B Evaluation Domains	155
References	159