

Time-Dependence in Markovian Decision Processes

Jeremy James McMahon

Thesis submitted for the degree of

Doctor of Philosophy

in

Applied Mathematics

at

The University of Adelaide

(Faculty of Mathematical and Computer Sciences)

School of Applied Mathematics



September, 2008

This work contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text.

I consent to this copy of my thesis, when deposited in the University Library, being made available in all forms of media, now or hereafter known.

SIGNED: DATE:

I extend sincere gratitude to my two supervisors, Professor Nigel Bean and Professor Michael Rumsewicz, who have both inspired and guided me to complete this work. Their input and encouragement has been extremely beneficial and I thank them wholeheartedly for their friendship and support over the last few years.

I am grateful to all three of my parents for the love and reassurance they have given me, not just throughout my time as a PhD student. Without them, any hope of reaching this point in my education would at best be a distant dream. I particularly appreciate the editorial contributions of my mother who, despite supposedly limited mathematical knowledge, provided invaluable feedback.

Lastly, I thank my beautiful wife who has endured much whilst I have been working on this thesis. I love and cherish Sarah for being by my side, encouraging me and always believing in me. Now we may begin the next chapter of our lives.

For Olive.

Contents

Abstract	xv
1 Introduction	1
2 Markov Processes	7
2.1 Introduction	7
2.2 A Markov Process	8
2.2.1 The Markovian Assumption	9
2.2.2 Time-Homogeneity	10
2.2.3 Analysis of Discrete-Time Markov Processes	10
2.2.4 Analysis of Continuous-Time Markov Processes	11
2.2.5 Discretizing Via Uniformization	14
2.2.6 Applications	16
2.3 A Markov Decision Process	19
2.3.1 Rewards and Decisions	20
2.3.2 Finite Horizon	22
2.3.3 Infinite Horizon	25
2.3.4 Continuous Time	27
2.4 A Semi-Markov Decision Process	30
2.5 A Generalized Semi-Markov Decision Process	33
3 Phase-Type Distributions	37
3.1 Introduction	37
3.2 Phase-Type Representations	38

3.3	Using Phase-Type Distributions	42
4	The Race	45
4.1	Introduction	45
4.2	The Race – Formal Description	47
4.3	Restricted Vision	49
4.3.1	Blind	49
4.3.2	Partially Observable	53
4.4	Full Vision	59
4.4.1	Value Equations	60
4.4.2	Policy Evaluation	62
4.4.3	The Race Revisited	64
4.5	The Race – Exponential System	67
4.5.1	Value Equations	68
4.5.2	MDP Approach	70
5	The Race – Erlang System	75
5.1	Introduction	75
5.2	Value Equations	77
5.2.1	State K	78
5.2.2	State $K - 1$	79
5.2.3	State $K - 2$	87
5.3	Summary	91
6	Phase-Space Model – Erlang System	95
6.1	Introduction	95
6.2	Existing Phase-Space Techniques	100
6.3	Our Phase-Space Technique	108
6.3.1	Level K	111
6.3.2	Level $K - 1$	111
6.3.3	Level $K - 2$	115

6.3.4	Level $K - 3$	123
6.4	Summary	127
7	Phase-Space Model – General Analysis	131
7.1	The Decision Process and Optimal Actions	131
7.2	Phase-Space Construction	133
7.3	Action-Consistent Valuation	140
7.4	Optimality Equations	144
7.5	Level-skipping in the Phase-Space	148
7.6	The Phase-Space Technique	154
8	Time-Inhomogeneous MDPs	157
8.1	Introduction	157
8.2	Time-Inhomogeneous Discounting	159
8.3	The Random Time Clock Technique	162
8.3.1	Time Representation	163
8.3.2	State-Space Construction	165
8.3.3	Reward Structure and Discounting	168
8.3.4	Truncation	171
8.3.5	Implementation	174
8.3.6	Extension for Time-Inhomogeneous Transitions	178
8.4	The Race – Erlang System	180
8.5	Summary	192
9	Conclusions	195
	References	199

List of Figures

2.2.1 Example of a continuous-time Markov process	17
2.5.1 State-space of the toast and tea example	35
3.1.1 Graphical representation of a selection of PH distributions	39
4.3.1 Optimal waiting time in state 2 given decision epoch at time s	52
4.3.2 Optimal expected value for state 2 at decision epoch s	57
4.3.3 Optimal expected value for state 1 at decision epoch s	58
4.3.4 Optimal expected value for state 0 at decision epoch s	59
4.5.1 Markov chain state-space of the exponential system	71
5.2.1 Expected value for state 2 at decision epoch s for differing actions	87
5.2.2 Optimal expected value for state 1 at decision epoch s	90
5.3.1 Optimal expected value for state 0 at decision epoch s	92
6.1.1 Markov chain representation of the Erlang order p distribution	96
6.1.2 Markov Chain representation of the K Erlang order 2 phase-space	98
6.2.1 Comparison of expected values for optimal and randomized policies	102
6.2.2 Guideline summary of the phase tracking model	103
6.2.3 Comparison of techniques for state/level 1	106
6.2.4 Expected optimal value of level 2 as seen from level 1	107
6.3.1 Guideline summary of the phase-space technique	110
6.3.2 Continuation values in level 1	121
6.3.3 Comparison of techniques for state/level 1	122
6.3.4 Comparison of techniques for state/level 0	126

6.4.1	Algorithmic summary of the phase-space technique for the race . . .	129
7.2.1	A two state semi-Markov reward process	135
7.2.2	Phase-space of a two state semi-Markov reward process	136
7.5.1	Level-skipping of a (TD, NAC) level	151
7.5.2	Level-skipping of a (TD, AC) level	151
7.5.3	Example of level-skipping in the phase-space technique	154
8.2.1	MOS decay as end-to-end delay is increased	160
8.3.1	State-space of a simple 2 state Markov process	163
8.3.2	Erlang density function of mean 1 with differing parameters	165
8.3.3	RTC State-space of a 2 state Markov process	167
8.3.4	Algorithmic summary of the RTC technique	177
8.3.5	RTC State-space of a 2 state time-inhomogeneous Markov process . .	179
8.4.1	Markov chain defined by $\mathbf{Q}_0(t)$	183
8.4.2	Markov chain defined by $\mathbf{Q}_1(t)$	183
8.4.3	Sigmoid absolute discount function	184
8.4.4	RTC state-space and transition rates when <i>continue</i> is selected	186
8.4.5	RTC technique with various time-state resolutions	188
8.4.6	Technique comparison for optimal value of state 2	190
8.4.7	Absolute error of the RTC technique for state 2	191
8.4.8	Technique comparison for optimal value of state 1	192

List of Tables

5.2.1 Summary of optimal policies	85
6.1.1 Optimal values for phase-states in a K Erlang order 2 system	99
8.4.1 Termination rewards for the RTC state-space	187

Abstract

The main focus of this thesis is Markovian decision processes with an emphasis on incorporating time-dependence into the system dynamics. When considering such decision processes, we provide value equations that apply to a large range of classes of Markovian decision processes, including Markov decision processes (MDPs) and semi-Markov decision processes (SMDPs), time-homogeneous or otherwise. We then formulate a simple decision process with exponential state transitions and solve this decision process using two separate techniques. The first technique solves the value equations directly, and the second utilizes an existing continuous-time MDP solution technique.

To incorporate time-dependence into the transition dynamics of the process, we examine a particular decision process with state transitions determined by the Erlang distribution. Although this process is originally classed as a generalized semi-Markov decision process, we re-define it as a time-inhomogeneous SMDP. We show that even for a simply stated process with desirable state-space properties, the complexity of the value equations becomes so substantial that useful analytic expressions for the optimal solutions for all states of the process are unattainable.

We develop a new technique, utilizing phase-type (*PH*) distributions, in an effort to address these complexity issues. By using *PH* representations, we construct a new state-space for the process, referred to as the phase-space, incorporating the phases of the state transition probability distributions. In performing this step, we effectively model the original process as a continuous-time MDP. The information available in this system is, however, richer than that of the original system. In the interest of maintaining the physical characteristics of the original system, we define

a new valuation technique for the phase-space that shields some of this information from the decision maker. Using the process of phase-space construction and our valuation technique, we define an original system of value equations for this phase-space that are equivalent to those for the general Markovian decision processes mentioned earlier. An example of our own phase-space technique is given for the aforementioned Erlang decision process and we identify certain characteristics of the optimal solution such that, when applicable, the implementation of our phase-space technique is greatly simplified.

These newly defined value equations for the phase-space are potentially as complex to solve as those defined for the original model. Restricting our focus to systems with acyclic state-spaces though, we describe a top-down approach to solution of the phase-space value equations for more general processes than those considered thus far. Again, we identify characteristics of the optimal solution to look for when implementing this technique and provide simplifications of the value equations where these characteristics are present. We note, however, that it is almost impossible to determine *a priori* the class of processes for which the simplifications outlined in our phase-space technique will be applicable. Nevertheless, we do no worse in terms of complexity by utilizing our phase-space technique, and leave open the opportunity to simplify the solution process if an appropriate situation arises.

The phase-space technique can handle time-dependence in the state transition probabilities, but is insufficient for any process with time-dependent reward structures or discounting. To address such decision processes, we define an approximation technique for the solution of the class of infinite horizon decision processes whose state transitions and reward structures are described with reference to a single global clock. This technique discretizes time into exponentially distributed length intervals and incorporates this absolute time information into the state-space. For processes where the state-transitions are not exponentially distributed, we use the hazard rates of the transition probability distributions evaluated at the discrete time points to model the transition dynamics of the system. We provide a suitable reward structure approximation using our discrete time points and guidelines for sensible truncation,

using an MDP approximation to the tail behaviour of the original infinite horizon process. The result is a finite-state time-homogeneous MDP approximation to the original process and this MDP may be solved using standard existing solution techniques. The approximate solution to the original process can then be inferred from the solution to our MDP approximation.