

A Bayesian Framework for Online Parameter Learning in POMDPs

Amin Atrash

Doctor of Philosophy

School of Computer Science

McGill University

Montreal, Quebec

May 20, 2011

A thesis submitted to McGill University in partial fulfillment of the requirements of
the degree of a Doctor of Philosophy

©Amin Atrash, 2011

ACKNOWLEDGEMENTS

Many people have helped me reach the completion of this thesis. I would like to thank my advisor, Joelle Pineau, for many years of advice and guidance. Her support and contributions were invaluable to my progress. I would also like to thank Doina Precup, Prakash Panangaden, and Gregory Dudek. It was an honor and privilege to work with these amazing people.

I would like to thank the members of the wheelchair team, Robert Kaplow, Julien Villemure, and Robert West, for their help and support with the project. The Reasoning and Learning Lab has been my academic home for many years, and I thank members such as Cosmin Padaruru, Pablo Castro, Jordan Frank, Mahdi Fard, Robert Vincent, Gheorge Comanici, Sophia Knight, Kamal Al-Marhoobi, Shaowei Png, Keith Bush, Arthur Guez, Monica Dinculescu, Stephane Ross, Caitlin Phillips, Hiba Yamani, Philip Bachman, Jonathon Taylor, Phillip Keller, Ryan Faulkner, and Marc Bellemare who supported me along the way. In addition, people such as Patrick and Joe Mattar, Mia Zentari, Ladan Mahabadi, Erin Dickie, Ilinca Popovici, Jean-Michel Petit, and Neil Barnaby helped make my time at McGill exciting. My thanks to Khai Troung whose experience I have drawn upon many times.

Finally, I would like to give special thanks to my parents, Hani and Nouhad Atrash, my brother and sister, Amer and Zeina Atrash, and all my relatives. Without their support, I would not be where I am today.

ABSTRACT

Decision-making under uncertainty has become critical as autonomous and semi-autonomous agents become more ubiquitous in our society. These agents must deal with uncertainty and ambiguity from the environment and still perform desired tasks robustly. Partially observable Markov decision processes (POMDPs) provide a principled mathematical framework for modelling agents operating in such an environment. These models are able to capture the uncertainty from noisy sensors, inaccurate actuators, and perform decision-making in light of the agent’s incomplete knowledge of the world. POMDPs have been applied successfully in domains ranging from robotics to dialogue management to medical systems.

Extensive research has been conducted on methods for optimizing policies for POMDPs. However, these methods typically assume a model of the environment is known. This thesis presents a Bayesian reinforcement learning framework for learning POMDP parameters during execution. This framework takes advantage of agents which work alongside an operator who can provide optimal policy information to help direct the learning. By using Bayesian reinforcement learning, the agent can perform learning concurrently with execution, incorporate incoming data immediately, and take advantage of prior knowledge of the world. By using such a framework, an agent is able to adapt its policy to that of the operator.

This framework is validated on data collected from the interaction manager of an autonomous wheelchair. The interaction manager acts as an intelligent interface between the user and the robot, allowing the user to issue high-level commands

through natural interface such as speech. This interaction manager is controlled using a POMDP and acts as a rich scenario for learning in which the agent must adjust to the needs of the user over time.

ABRÉGÉ

Comme le nombre d'agents autonomes et semi-autonomes dans notre société ne cesse de croître, les prises de décisions sous incertitude constituent désormais un problème critique. Malgré l'incertitude et l'ambiguïté inhérentes à leurs environnements, ces agents doivent demeurer robustes dans l'exécution de leurs tâches. Les processus de décision markoviens partiellement observables (POMDP) offrent un cadre mathématique permettant la modélisation des agents et de leurs environnements. Ces modèles sont capables de capturer l'incertitude due aux perturbations dans les capteurs ainsi qu'aux actionneurs imprécis. Ils permettent conséquemment une prise de décision tenant compte des connaissances imparfaites des agents. À ce jour, les POMDP ont été utilisés avec succès dans un éventail de domaines, allant de la robotique à la gestion de dialogue, en passant par la médecine.

Plusieurs travaux de recherche se sont penchés sur des méthodes visant à optimiser les POMDP. Cependant, ces méthodes requièrent habituellement un modèle environnemental préalablement connu. Dans ce mémoire, une méthode bayésienne d'apprentissage par renforcement est présentée, avec laquelle il est possible d'apprendre les paramètres du modèle POMDP pendant l'exécution. Cette méthode tire avantage d'une coopération avec un opérateur capable de guider l'apprentissage en divulguant certaines données optimales. Avec l'aide du renforcement bayésien, l'agent peut apprendre pendant l'exécution, incorporer immédiatement les données nouvelles et profiter des connaissances précédentes, pour finalement pouvoir adapter sa politique de décision à celle de l'opérateur.

La méthodologie décrite est validée à l'aide de données produites par le gestionnaire d'interactions d'une chaise roulante autonome. Ce gestionnaire prend la forme d'une interface intelligente entre le robot et l'utilisateur, permettant à celui-ci de stipuler des commandes de haut niveau de façon naturelle, par exemple en parlant à voix haute. Les fonctions du gestionnaire sont accomplies à l'aide d'un POMDP et constituent un scénario d'apprentissage idéal, dans lequel l'agent doit s'ajuster progressivement aux besoins de l'utilisateur.

TABLE OF CONTENTS

| | |
|--|-----|
| ACKNOWLEDGEMENTS | ii |
| ABSTRACT | iii |
| ABRÉGÉ | v |
| LIST OF TABLES | x |
| LIST OF FIGURES | xi |
| 1 Introduction | 2 |
| 1.1 Application Domain | 4 |
| 1.2 Probabilistic Models for Decision-Making | 5 |
| 1.3 Learning | 7 |
| 1.4 Bayesian Reinforcement Learning | 8 |
| 1.5 Learning From Examples | 10 |
| 1.6 Contributions | 10 |
| 1.7 Organization | 12 |
| 2 Smartwheeler Autonomous Wheelchair | 13 |
| 2.1 Overview of the Smartwheeler Autonomous Wheelchair | 14 |
| 2.2 Platforms | 16 |
| 2.3 Interaction Interface | 18 |
| 2.3.1 Speech | 19 |
| 2.3.2 Semantic Grammar | 20 |
| 2.3.3 Interaction Manager | 21 |
| 2.3.4 Examples | 22 |
| 2.4 Target Population | 24 |
| 2.5 Wheelchair Skills Test | 25 |
| 2.6 Performance | 26 |
| 2.7 Discussion | 28 |
| 2.8 Related Work | 29 |

| | | |
|-------|---|-----|
| 2.8.1 | Autonomous Wheelchairs | 29 |
| 2.8.2 | Dialogue Management | 36 |
| 2.9 | Future Work | 41 |
| 3 | Probabilistic Models for Sequential Decision-Making | 42 |
| 3.1 | Markov Decision Processes | 42 |
| 3.1.1 | Value Iteration | 45 |
| 3.2 | Partially Observable Markov Decision Processes | 46 |
| 3.2.1 | Belief State Tracking | 47 |
| 3.2.2 | Exact Value Iteration | 49 |
| 3.2.3 | Point-based Approximate Methods | 54 |
| 3.3 | Reinforcement Learning | 57 |
| 3.3.1 | Reinforcement Learning in Fully Observable Domains | 58 |
| 3.3.2 | Learning in Partially Observable Domains | 59 |
| 3.4 | Summary | 63 |
| 4 | Bayesian Reinforcement Learning | 64 |
| 4.1 | Overview | 66 |
| 4.2 | Bayesian Reinforcement Learning in Fully Observable Domains | 69 |
| 4.3 | Bayesian Reinforcement Learning in Partially Observable Domains | 72 |
| 4.4 | Summary | 74 |
| 5 | A Framework for Online Parameter Learning | 76 |
| 5.1 | Framework | 76 |
| 5.2 | Policy Oracle | 80 |
| 5.3 | Representation | 81 |
| 5.4 | Distribution Update | 83 |
| 5.5 | Action Selection | 85 |
| 5.6 | Example | 86 |
| 5.7 | Discussion | 87 |
| 6 | Learning Reward | 89 |
| 6.1 | Representing and Learning Reward | 90 |
| 6.2 | Experimental Methodology | 91 |
| 6.3 | Results | 95 |
| 6.4 | Related Work | 107 |
| 6.5 | Conclusion | 108 |

| | | |
|-----|--|-----|
| 7 | Learning Observations | 112 |
| 7.1 | Representing and Learning Observations | 113 |
| 7.2 | Learning Rewards and Observations | 114 |
| 7.3 | Evaluation | 115 |
| 7.4 | Related Work | 125 |
| 7.5 | Future Work | 126 |
| 7.6 | Conclusion | 126 |
| 8 | Conclusion | 128 |
| 8.1 | Framework | 129 |
| 8.2 | Experimental Results | 130 |
| 8.3 | Limitations and Future Work | 131 |
| 8.4 | Summary | 134 |
| | Appendix A - Distributions | 136 |
| | Appendix B - Interaction Manager POMDP | 138 |

LIST OF TABLES

| <u>Table</u> | | <u>page</u> |
|--------------|---|-------------|
| 2-1 | Semantic Slots and Possible Value Assignments | 20 |
| 2-2 | Results for healthy subjects | 35 |
| 2-3 | Results for target subjects | 35 |

LIST OF FIGURES

| <u>Figure</u> | <u>page</u> |
|--|-------------|
| 2-1 Autonomous Wheelchairs | 16 |
| 2-2 Architecture of the Smartwheeler Autonomous Wheelchair | 18 |
| 2-3 Example of dialogue involving a query | 23 |
| 2-4 Example of dialogue with misrecognition | 24 |
| 3-1 Example Markov Decision Process | 43 |
| 3-2 Example Value Function | 50 |
| 6-1 Learning Reward on Interaction Manager - <i>Data-based</i> | 96 |
| 6-2 Learning Reward on Interaction Manager - <i>Generative-based</i> | 97 |
| 6-3 Varying Number of Models Resampled - <i>Data-based</i> | 99 |
| 6-4 Varying Number of Models Resampled - <i>Generative-based</i> | 100 |
| 6-5 Results using only resampling | 102 |
| 6-6 Results using all models to update the distribution | 103 |
| 6-7 Impact of priors on learning | 104 |
| 6-8 Results using supervised learning | 105 |
| 7-1 Learning Observations | 118 |
| 7-2 Comparison against Baum-Welch | 119 |
| 7-3 Learning Observations with Error in Rewards | 120 |
| 7-4 Learning Rewards and Observations | 122 |
| 7-5 Learning over SACTI Domain | 123 |

Notation

- S : set of states
- s_t : state at time t
- s_i : i^{th} state of the model
- A : set of actions
- a_t : action at time t
- a_i : i^{th} action of the model
- Z : set of observations
- z_t : observation at time t
- z_i : i^{th} observation of the model
- π : policy
- r_t : reward received at time t
- γ : discount factor
- b : belief

- $T()$: transition function
- $R()$: reward function
- $O()$: observation function

- $V()$: value function for state s
- $V^t()$: value function for state s at time t
- $Q()$: value of execution action a in state s
- α : vector defining hyper-plane for POMDP value function
- Γ : set of α -vectors

- M_i : i^{th} sampled model

CHAPTER 1

Introduction

Personal robots have become more ubiquitous in recent years and will continue to do so as technology increases their capabilities and opportunities for deployment. As robots move from industrial settings to everyday environments, they face many new challenges. Robots in industrial settings work under very controlled conditions and perform a limited set of repetitive tasks. These robots are often controlled by a well trained operator. Personal robots are expected to handle a more dynamic and uncontrolled environment essentially designed for humans, and must interact with users who have little or no training. These robots must learn about the environment, plan to handle tasks, and adapt to changes. This is challenging due to a degree of uncertainty in the environment, from sensor noise and ambiguity, to unpredictable actions.

Personal robots face further challenges as they are required to interact with people on a regular basis. The interaction should be as natural as possible to minimize the burden on the users. Humans are accustomed to interacting with each other through speech and gestures. Other natural interfaces can also be considered, such as touch screen displays and other intuitive devices. Robots need to be deployed in situations where technical expertise for controlling the robot may not be available. Such scenarios may include, for example, a disabled user with a robotic wheelchair, or a diver accompanied by an autonomous underwater robot, where interaction may

be limited. Similarly, a user attempting to manage multiple agents will need to distribute his attention among many sources. A higher degree of autonomy in each agent, combined with an intuitive interface, would aid the operator of the system.

One of the core challenges in improving robot autonomy is the ability for an agent to learn and adapt to its environment. Many decision-making techniques assume a model of the environment is provided. However, as modern domains become more complex, they become increasingly difficult to model by hand. Instead, it is advantageous for an agent to develop its own model of the environment. This reduces the need for a human to input the information in advance. This is particularly important in domains which may be too complex for a human to specify. Additionally, this allows the agent to operate in dynamic environments, whereas a static model provided to the robot may become outdated as the environment changes.

In this thesis, we address the issue of learning in partially observable domains. We present the use of Bayesian reinforcement learning to represent the uncertainty over unknown parameters. We propose an approach which utilizes oracle information to guide the learning, and address the issue of action selection given such a representation. We also present work on an autonomous wheelchair, detailing the interface design as well as user studies. These studies act as motivation as well as providing a data set for validation. We will demonstrate the ability of the framework to determine the handcrafted model automatically using the data collected from user studies.

1.1 Application Domain

Consider the Smartwheeler autonomous wheelchair (Pineau & et. al., 2007). This project focuses on the development of a robotic wheelchair for use by a disabled population by taking advantage of advancements in robotics and artificial intelligence. Here, the goal is to reduce the physical and cognitive load on patients with mobility limitations by shifting partial control to an imbedded intelligent controller.

Clearly in this case, requiring a sophisticated control input system is unreasonable. Ideally, the user should be able to interact and control his wheelchair as naturally as he would interact with a human assistant moving him through a building. The wheelchair should be expected to respond to traditional communication modalities used between humans, such as speech and gestures, learn about the specific environment in which it is deployed, and adapt to the user's preferences over time. The wheelchair should be able to handle this without unnecessarily burdening the user.

The interface of this wheelchair must address the difficulties of interaction with a human operator. This interface must act as a gateway between the user and the robot control system by receiving multi-modal user input and determining the proper action to relay to the robot or response system. This interface must handle a large amount of uncertainty from the world. Speech is an obvious natural interface, but speech recognition systems are imperfect and may not accurately capture what is said by the user. Even if recognition is perfect, commands to the wheelchair can contain ambiguity or incomplete information. Similarly, the user's commands cannot be known ahead of time. However, accounting for the possible requests is crucially

important for decision making. The system must consider all these factors when attempting to determine the user's needs, while at the same time selecting the best actions.

Finally, the wheelchair interface must be able to adapt to each individual user. The preference of each user operating such a wheelchair will vary between individuals. Handcrafting a model on a per user basis would be a lengthy and time consuming task. Instead, the wheelchair control system should adapt its own parameters for new users as they use the wheelchair. In a similar fashion, the control system should adapt to users whose preferences change over time. Prolonged use of the wheelchair may result in different expectations in early stages versus the long term, as users become more comfortable with the chair, or simply change their own requirements.

1.2 Probabilistic Models for Decision-Making

Markov Decision Processes (Bellman, 1957) are probabilistic models for decision making in domains with uncertainty. These types of models are able to perform decision making in light of uncertainty from the environment. Markov Decision Processes define the world as a set of *states*. An agent transitions nondeterministically between states by executing *actions*. These actions are modelled using *transition* functions which capture the uncertainty in the outcome of each action. Once an action is executed, the agent receives a *reward* based on the action taken and resulting state. Rewards represent the cost of taking an action, such as the time required to execute the action or reward received for reaching a goal state. By using this framework to model the environment, an agent can perform decision-making in light of the uncertainty from the environment.

This thesis focuses on a general form of Markov models known as Partially Observable Markov Decision Processes (POMDPs) (Sondik, 1971; Smallwood & Sondik, 1973). POMDPs generalize the MDP framework by assuming the environment is *partially observable*. In the case of MDPs, the world is assumed to be *fully observable*, in that the agent is aware at all times of the current state of the environment. POMDPs remove this assumption and instead rely on noisy sensory information, known as *observations*, to infer the underlying state. These observations attempt to model the various sensors of the robot. Inferring the underlying state is typically accomplished by maintaining a probability distribution over the potential states of the world known as the *belief state*. By combining the actions taken and the observations sensed, the belief state of the agent can be updated and maintained. This approach is often a more realistic representation of problems and can account for more source of uncertainty than the fully observable case, providing an additional level of robustness. Due to this robustness, POMDPs have been shown useful for representation and decision-making in a range of applications (Cassandra, 1998), including robot navigation (Koenig & Simmons, 1996b), dialogue management (Williams & Young, 2005; Pineau et al., 2003), and behavior tracking (Hoey et al., 2005).

Consider the use of POMDPs to model the Smartwheeler interface described in Section 1.1. In this case, the states can capture the user’s true intention, what the user wishes the robot to execute. However, the speech recognition system provides noisy recognition, treated as observations. The interface must use this information to select the best action. Here, actions are commands issued to the robot, or possibly requests by the interface for additional information from the user when uncertainty

is prevalent. A policy over this model dictates when the interface should issue a command to the robot, or request clarification from the user.

1.3 Learning

While a variety of algorithms exist for decision-making using POMDPs (Sondik, 1971; Pineau et al., 2003), learning the models themselves has not been as thoroughly examined. Decision-makers typically assume a complete model of the world is provided. In real-world applications, providing this detailed level of information is not realistic. The domains may be extensive and complex, and attaining proper sensor models may be costly. Proper models may require input and representation of expert knowledge, or the collection of considerable data sets. A planner may only have access to a rough model of the world as provided by a user, and in the worst case, there may be no information at all. To be effective in real-world domains, an agent must be able to accurately estimate and learn the parameters of the POMDP representation of the world through experience.

Initial research in POMDP learning drew on work from the closely related Hidden Markov Models (HMMs) and focused on application of the Baum-Welch algorithm (Rabiner, 1990). Hidden Markov Models are similar to POMDPs, but are only capable of recognition of a situation, and are unable to provide decision-making. Baum-Welch uses an expectation-maximization approach to estimate the parameters of the model. Early work (Chrisman, 1992) demonstrated the use of Baum-Welch, traditionally applied to HMMs, to learn the parameters of POMDPs. Similar approaches have been explored for learning models during robot navigation (Koenig & Simmons, 1996b; Shatkey & Kaelbling, 1997; Theodorou et al., 2001; Nikovski

& Nourbakhsh, 2000; Kurniawati et al., 2009; Song et al., 2010). However, because HMMs only focus on recognition, these approaches are only able to estimate the observation and transition parameters of the models, leaving the reward parameters to be specified by hand. Modern learning algorithms attempt to leverage policy information to determine an appropriate reward function through the use of *reinforcement learning* (Sutton, 1988; Bertsekas & Tsitsiklis, 1996; Moore & Atkeson, 1993; Sutton & Barto, 1998). However, these techniques are not practical in partially observable domains.

While probabilistic models are an ideal representation for the interface of the Smartwheeler, constructing an accurate model by hand is difficult for a such a complex domain. A designer must consider the user’s behavior, models of the sensors, and the positive and negative effects of each action. Even after a model is created, there are no assurances of its applicability to more than a single user. Thus, the ability for an agent to learn and adapt is essential to allow such a system to be robust. By taking advantage of learning techniques, this information could be determined automatically by the agent. Reinforcement learning would allow the agent to adjust the model parameters as it interacts with the user. This would also provide a mechanism to adjust the model to different users. A designer may begin with an initial crude representation, which is then fine-tuned automatically when deployed.

1.4 Bayesian Reinforcement Learning

Traditional HMM and reinforcement learning methods maintain the learned values as point estimates. This fails to capture the quality of the values learned. For example, differentiating between values learned by an agent which has just begun

operating versus an agent which has had significant experience in the domain would be difficult. Conventional approaches apply an external mechanism (Kaelbling et al., 1996) to help select when an agent should execute an action to explore the environment to gain further knowledge, or when to select the action the agent currently believes is best. This is known as the exploration versus exploitation problem.

Recent work has sought to integrate Bayesian reasoning with reinforcement learning (Dearden et al., 1998; Dearden et al., 1999; Duff, 2002). This is known as *Bayesian reinforcement learning*. Bayesian reinforcement learning seeks to represent and maintain the parameters as distributions. This provides several advantages. By maintaining distributions over the parameters, an agent is able to reason about the uncertainty in the known models, for example, considering the variance in the potential values of a parameter when making a decision. This approach also allows the integration of prior knowledge. Distributions can be initialized using knowledge of the environment, allowing the system, for example, to begin with a rough model of the world.

In the case of the interface for the autonomous wheelchair, modeling the details of the interface may require extensive domain knowledge and repeated trials to adjust the behavior. Using a Bayesian reinforcement learning framework, the developer can implement a rough model which is then refined through experience. This model can furthermore be adjusted differently for each wheelchair on which it is deployed.

While a Bayesian approach to learning provides many benefits, several issues are raised. First, the choice of distributions to represent the parameters plays a critical role in the learning. Second, the distributions must be updated based on experience.

Finally, the agent must be able to select an action based on the distributions. While there exists an extensive literature examining the use of Bayesian reinforcement learning towards the case of fully observable MDPs, its application towards POMDPs has been less explored. One of the primary goals of this thesis is to further the work on Bayesian reinforcement learning using POMDPs, and explore issues such as representation, action selection, and model updates in this context.

1.5 Learning From Examples

While agents using reinforcement learning are capable of learning independently, this process can require extensive data and numerous attempts. These trials may be expensive, time consuming, or even dangerous, particularly in the case of robotics. Consider examples of robots working alongside humans, as in the wheelchair domain described above, where the human element increases the complexity of the problem, as well as requiring the robot to adapt very quickly to the needs of the user. Although agents working alongside humans face many additional challenges, this situation also provides an opportunity to gain useful information to aid with learning. An operator can provide information to the agent to help direct the learning process. The framework in this thesis will take advantage of this operator in the form of an *oracle*. The goal of the framework is to adjust the parameters to result in behavior similar to that of the oracle. This framework will demonstrate a mechanism to utilize this information to help guide the learning.

1.6 Contributions

The major contributions of this thesis are as follows:

- **A unified framework for learning POMDP parameters** - We present a unified framework for learning POMDP parameters based on Bayesian reinforcement learning. Earlier work on BRL for POMDPs has centered around learning a single class of parameters, either rewards, transitions, or observations. The case of learning the reward parameters is particularly difficult. The framework presented in this thesis allows for a unified approach for learning multiple types of parameters concurrently, including the reward parameters. This framework utilizes a sampling approach to approximate the distributions. By sampling, the framework can remain efficient, and the samples can be utilized by the learning mechanisms.
- **Use of policy oracles** - We describe the use of an operator acting as an oracle to obtain information about the environment. We assume the agent has access to an oracle which can provide the agent with the correct policy information. We present a novel approach to take advantage of the policy information provided to direct the model learning.
- **Updating model parameters** - We propose a new technique to update the distribution parameters during learning. This method combines an approximation of the distributions with the information from the oracle to update the models.
- **Validation on user data collected from the Smartwheeler Autonomous Wheelchair** - A major emphasis of this thesis lies in the validation of the framework in a real-world domain. While other work has examined the use of BRL or learning in POMDPs, experiments are limited to simple example

domains or small real world domains. We demonstrate our framework on data collected during user studies of the Smartwheeler Autonomous Wheelchair.

This thesis represents the culmination of work in the development of this framework. Earlier work has demonstrated the use of the framework for learning only rewards (Atrash & Pineau, 2009). Validation of this earlier work was limited to smaller example POMDP domains as well as a simplified dialogue domain. This thesis continues this work by demonstrating learning of multiple parameters in unison, as well as validation on a more complex domain. Preliminary work on the Smartwheeler demonstrated an early version of the interface on a small population of students (Atrash et al., 2009). This thesis presents recent user studies while describing the autonomous wheelchair domain.

1.7 Organization

The remainder of this thesis is laid out as follows: The Smartwheeler autonomous wheelchair domain is detailed in Chapter 2. This includes details of the architecture as well as a description of the user studies conducted. In Chapter 3, we review MDPs, POMDPs, and the associated planning algorithms. Chapter 4 presents a survey of existing Bayesian reinforcement learning approaches. A Bayesian framework for learning is presented in Chapter 5. Here, we discuss the issues associated with the framework including representation, action selection, and updates of the distributions. Experimental validation of the framework appears in Chapters 6 and 7. Finally, we conclude the discussion of this work in Chapter 8.

CHAPTER 2

Smartwheeler Autonomous Wheelchair

In this chapter, we present an overview of the Smartwheeler autonomous wheelchair project (Atrash et al., 2009). The goal of the Smartwheeler project is to develop an autonomous wheelchair to aid individuals with mobility impairments. Controlling a robot such as this requires a unique interface, in this case an intelligent interaction interface. This interaction interface acts as the layer between the user and the lower level robot control, taking sensory information and determining an appropriate action.

This interface provides a rich environment for demonstrating the learning framework on a deployed system. Decision-making must be conducted in light of noise from the input as well as uncertainty about the user’s true intention. These characteristics lend themselves towards POMDPs as a viable representation. POMDPs have previously been demonstrated as a common tool for the closely-related dialogue management domain. Models for such a complex task would be difficult to create by hand, and may vary from user to user. Hence, learning the models is an important challenge for such a system. An agent controlling an autonomous wheelchair naturally operates in an environment alongside a human operator who can provide information to the agent when needed.

The purpose of this chapter is to present an overview of the interaction system, describe the experimental procedure used for evaluation, and discuss results of the

performance of the system. The data collected is used for validation of the learning framework in later chapters.

2.1 Overview of the Smartwheeler Autonomous Wheelchair

Many individuals with mobility limitations require the use of an electric wheelchair to navigate through their environment, with as many as 4.3 million power chair users in the United States alone (Simpson et al., 2008). Recent surveys (Fehr et al., 2000) suggest that as many as 10% of these individuals experience difficulty maneuvering the power wheelchair. This is often a result of physical or cognitive fatigue, debilitating conditions, or sensory impairments that may limit the use of standard electric wheelchairs. It would be advantageous to these groups to reduce the burden of control.

The goal of the Smartwheeler autonomous wheelchair is to bring together developments in robotics and intelligent systems technology to reduce the physical and cognitive load on individuals with mobility impairments. Individuals with disabilities may find it physically exhausting to maneuver through cluttered rooms or tight doorways. Similarly, individuals with cognitive disabilities may be unable to respond quickly to rapidly changing environments such as crowded hallways. Recent developments in robotics technology have provided capabilities such as mapping and navigation, as well as development in sensor technology, which facilitate robots operating in everyday environments. By leveraging these developments, we can shift the burden of control from the user to the autonomous wheelchair. Instead of requiring a user to carefully position the wheelchair through a narrow door, or maneuver

through a congested hallway, these tasks can be shifted to the wheelchair’s intelligent controller.

Controlling such a sophisticated system requires a unique interface. Requiring the user to display computer expertise is clearly not a viable solution. Instead, the system should allow the user to issue high level tasks through common natural interfaces. Such a system poses many challenges. The interaction must be natural and comfortable for the user, requiring minimal technical expertise. The user should be able to communicate with the wheelchair as if dealing with a human assistant. The user should not be required to physically guide the wheelchair, instead offering high level goals and intentions which are then executed by the wheelchair itself. This leads to issues such as determining the user’s true intention in light of a very noisy sensing system, maintaining the context of the system over time, and making decisions given pieces of information from multiple sources.

In this chapter, we present the Smartwheeler autonomous wheelchair. This platform takes advantage of an intelligent multimodal interface which relies on simple modes of interaction, in this case, speech recognition and a touch screen display. Key to this interface is an underlying tracking and decision-making component. This component is able to incorporate the noise from the sensors, the uncertainty from the user’s intentions, and provide robots with decision-making capabilities. We will also detail the deployment and the evaluation of the autonomous wheelchair. An important element of this work is the evaluation of the intelligent interface on genuine user subjects. These subjects provide a realistic domain outside of a controlled research environment and represent the true target population of such a project. To



Figure 2–1: Autonomous Wheelchairs developed at McGill and Ecole Polytechnique de Montreal. These represent standard electric wheelchairs modified with onboard computing and sensing to allow for intelligent control.

this end, experiments were conducted in a rehabilitation center with an appropriate target population.

We will begin with a description of the target population and motivating evaluation domain. Both the population and domain are selected to be representative of a realistic domain. We will detail the Smartwheeler platform including the hardware and the intelligent interface design. Finally, we user study results describing the accuracy of the system when selecting actions with two groups of test subjects.

2.2 Platforms

Two similar robot wheelchair platforms were developed for this project, the first developed at McGill University and the second at École Polytechnique de Montréal. They can be seen in figure 2–1. The bases of the wheelchairs are a Sunrise Quickie Freestyle and a Sunrise Quickie 646 respectively. To allow for autonomous control, onboard computing, sensors, and actuators have been added to each wheelchair. Low power computers are mounted beneath the seats to provide on-board processing.

Laser measurement sensor (LMS) range finders are attached to the front and rear of each wheelchair to provide accurate distance readings around the wheelchair. Wheel encoders provide basic odometry during motion. Control is provided by CARMEN (Montemerlo et al., 2003) on the first wheelchair, and Player (Collett et al., 2005) on the second wheelchair. These robot control systems provide a set of algorithms for many well understood robot problems. Detailed maps can be constructed using the range-finders and the wheel encoders, which can then be used for localization and navigation. Obstacle avoidance is also implemented for local navigation.

To allow for a multimodal interface, devices were added to the wheelchairs to provide communication between the wheelchair and the user. A touch sensitive LCD screen is mounted on the right arm rest to display information to the user. A microphone is attached allowing the user to speak directly to the computer. Currently, this acts as the primary interface to the system. Finally, the original joystick of the wheelchair is functional and can be used to directly drive the wheelchair if necessary.

The current system focuses on the speech system as the primary interface, with the display providing feedback to the user when necessary. In the future, the multimodal system will be extended to take advantage of other input modes. Menus and maps can be displayed on the touchscreen display to be selected by the user. The joystick currently controls the wheelchair directly, but can be used as a selection device for menus on the display. A speech generation system could be added to allow the wheelchair to respond verbally to user. Finally, the rehabilitation community utilizes a large set of specialized input devices depending on a user's specific condition. These devices could also be integrated into the system based on a user's needs.

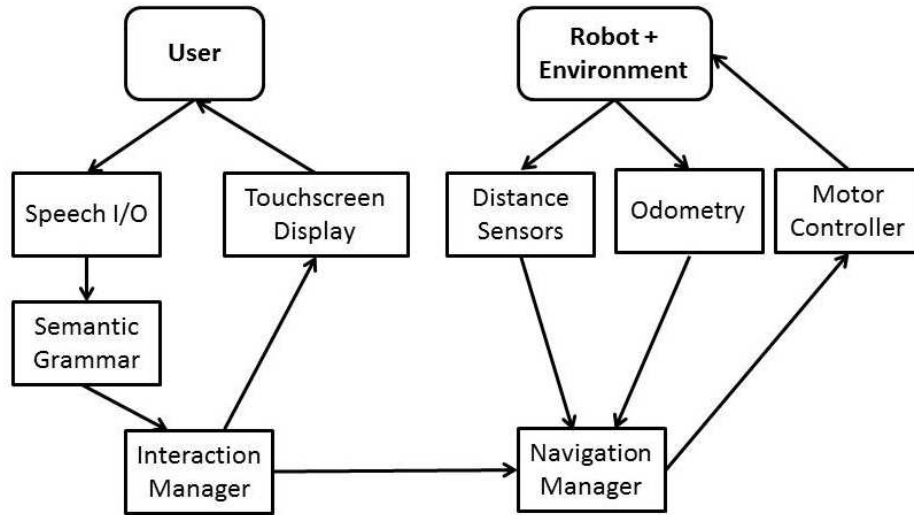


Figure 2–2: Architecture of the Smartwheeler Autonomous Wheelchair

2.3 Interaction Interface

Figure 2–2 presents an overview of the architecture for the intelligent interface of the Smartwheeler wheelchair. This architecture contains several components to process the input from the user. A user speaks to a speech recognition system which attempts to transcribe the audio. This transcription is passed to a semantic parser which attaches semantic and syntactic information if possible. The annotated information is then handled by the interaction manager, which acts as the decision-making component of the system. Feedback is provided to the user through the mounted display. Modules communicate using Acropolis (Zalzal et al., 2009), a framework for interprocess communication that allows modules to be easily replaced. Once an action is selected, a command is sent to the Navigation Manager which handles the execution of the action by the robot platform.

The purpose of this interface is to provide a method for robust control of the wheelchair for users who have difficulty operating the wheelchair manually. The design of this interface is critical to compensate for the noise that would otherwise appear in a simple speech driven system. Earlier work on test subjects in this domain (Atrash et al., 2009) demonstrated word misrecognition error rates of about 63%. The experiments conducted in this chapter reflect similar experiences. Misrecognitions are common, especially in the everyday environments the wheelchair may encounter. Users may have difficulty speaking due to their disabilities. Even in conditions where recognition is perfect, language contains ambiguity. For example, a phrase such as “turn the chair” may be recognized correctly, but would still lack information such as direction or speed to allow the wheelchair to respond correctly. The interface must account for all these sources of input while attempting to determine the user’s true intention.

In the following sections, we provide an overview of this interface.

2.3.1 Speech

The first step of the interface involves the speech recognition system. For this work, Dragon Naturally Speaking 10 (Nuance, 2009) was used. Dragon is a commercially available speech recognition package. As part of the interface, the user speaks a command, which Dragon attempts to transcribe into text. As noted before, this alone is not sufficient to command the wheelchair, due to misrecognition errors and ambiguities. Instead, the transcription is further analyzed by the remaining components.

| Semantic Slot | Possible Values |
|----------------|--|
| ActionValue | move, switch, activate, deactivate, select, set, drive, turn, NULL |
| ActionType | move-action, hardware-action, config-action, NULL |
| DirectionValue | away, back, forward, backward, right, left, NULL |
| DirectionType | direction, vector-direction, angle-direction, NULL |
| PatientType | controller, speed, drive-mode, seat, motor, NULL |
| PatientValue | wc-hardware, wc-property, NULL |

Table 2–1: Semantic Slots and Possible Value Assignments

2.3.2 Semantic Grammar

Once the speech is transcribed by the speech recognition system, the transcription is processed by the semantic grammar module (White, 2001). The aim of this module is to constrain the output of the speech recognizer by performing a semantic and syntactic analysis. This allows multiple phrases with the same meaning to be mapped to each other, as well as removing extraneous words from the phrase. For example, “move backward,” “drive backward,” and “let’s move the wheelchair backward” would all map to the same semantic representation. By performing this analysis, the complexity in the remaining stages can be reduced. Given a phrase, the semantic grammar outputs a vector of semantic slots with each slot taking on a value. Examples of semantic slots and their possible values can be seen in table 2–1, although a larger set is used in application. Typically, a state has 3 to 6 semantic slots assigned. The choice of slots and values varies based on the phrase.

For this work, a reduced grammar limited to the domain of the wheelchair operation, is used. If a phrase can be parsed, the annotated information is sent to the interaction manager. Otherwise, the phrase is treated as a *bag of words*. When

using a bag of words approach, each phrase is treated as an unordered group of words with no grammatical context. Each word is treated as an independent observation with no regards to the surrounding context. The bag of words is less desirable due to the reduced information, but may still contain useful information if the phrase cannot be parsed by the semantic grammar module.

2.3.3 Interaction Manager

The interaction manager acts as the decision-making component of Smartwheeler interface. This module takes the output of the semantic grammar module and aims to select the best action in response. To select the action, the interaction manager attempts to infer the user’s intent using noise models of the speech recognition system, models of the semantic grammar, and models of the user’s behavior. Based on the input, the interaction manager may decide to issue an appropriate action to the robot if the system is confident of the user’s intent, or else it can request additional information in the form of a query to the user. The interaction manager is able to select from 29 commands related to the wheelchair control. This set includes low-level control, such as “move forward” and “turn around” as well as high level commands, such as “go through the doorway” and “follow the right wall.” Currently, two forms of queries are implemented, a specific query which request parameters such as distance or direction, and a general query. Specific queries occur when there is some degree of ambiguity in the parameters of an action, such as understanding the requested command is a “move,” but lacking the direction. A general query calls for a repetition of the command. In this case the user can simply restate the command, or rephrase the command in a different fashion.

The interaction manager is modelled using a Partially Observable Markov Decision Process (POMDP). A technical presentation of POMDPs is provided in chapter 3.2. Details of the POMDP implementation of the Interaction Manager can be found in Appendix B. The parameters for the model were calculated using annotated transcriptions from previous work (Atrash et al., 2009), transcriptions collected during development of the system, and through hand-crafted integration of domain knowledge by the developers. This process involved initial collection of data beforehand, and repeated trials and modifications of the parameters until the model behaved robustly.

2.3.4 Examples

This section presents two examples of an interaction between the user and the intelligent interface. The examples present the flow and format of information during the interaction. The transcribed audio can be seen in lines labelled *AllPhrases*. The ground truth transcription as determined by the researchers is labelled as *Actual*. This information is not available to the system but is presented here for clarification. Phrases which are successfully parsed by the semantic grammar have *SentenceRecognized* set to *true*, followed by a series of additional semantic tokens and values. The action as selected by the interaction manager is presented as *Selected Action*.

Figure 2–3 demonstrates a case where the agent requests additional information via a specific query. In this case, the agent determines that “move backward” does not uniquely identify an intention, and thus requests move information. The user repeats “backward” which helps isolate the command. While at first appearance, it does not seem clarification is required, there are several commands which are variations of

```
Actual: "Move backward"
<SentenceRecognized value="true" />
<Phrase value="move backward" />
<AllPhrases value="[move backward]" />
<Action value="move" type="move-action" />
<SpeechAct value="command" />
<Direction value="backward" type="vector-direction" />
```

Selected Action: Move Query

```
Actual: "Backward"
<SentenceRecognized value="false" />
<Phrase value="backward" />
<AllPhrases value="[backward]" />
```

Selected Action: Move Backward

Figure 2-3: Example of dialogue involving a query

“move backward.” Based on the sensitivity of the system and the preferences of the user, the agent may or may not request the clarification. These are characteristics controlled by the parameters of POMDP. In this example, the system is being overly cautious.

This example also demonstrates the effect of the semantic grammar. Here, the first command was successfully parsed for information, allowing the system to determine the type of command as well as parameters. The second command, in this case a single word, was not parsable and thus treated as a *bag of words* by the interaction manager.

In this second example found in figure 2-4, the speech recognition system has significant error in the recognition. However, the agent is still able to determine the

```
Actual: "Let's move towards the wall and end up parallel to the wall"  
<SentenceRecognized value="false" />  
<Phrase value="" />  
<AllPhrases value="[the move to wall to wall and and and up  
veer overcome wall on]" />
```

Selected Action: Align to Wall

Figure 2-4: Example of dialogue with misrecognition

correct action. This is likely due to the word “wall” which, in our current system, is unique to the action *Align to wall*.

2.4 Target Population

To evaluate the Smartwheeler interface, a standardized target population was selected based on the users’ abilities, rather than their pathology. Users have varying degrees of proficiency operating a wheelchair. Others may have difficulty speaking due to their condition. The goal is to target a set of users who share the same characteristics, not simply those who require a wheelchair. This allows for a comparison of their performance using the system without a bias due to the conditions themselves. This is particularly important when conducting experiments on a small sample size of 8 to 10 subjects.

For this work, users were evaluated based on a set of criteria such as mobility impairments, postural limitations, and level of autonomy in daily life. Users were also expected to be fluent in English and able to speak clearly. This was a necessary constraint as some potential subjects had difficulty speaking due to their disability. Subjects were recruited through a rehabilitation center by occupational therapists seeking out potential subjects who met the specified criteria. Final eligibility was

determined based on a survey administered by a trained occupational therapist in advance of the tests.

2.5 Wheelchair Skills Test

To standardize the evaluation of our interface, we utilize a procedure developed in the rehabilitation community. The Wheelchair Skills Program (WSP)(Kirby et al., 2004; Kirby, 2007) was developed at Dalhousie University to assess the capabilities of users operating a wheelchair through a set of representative tasks. The goal of this program is to identify the skills to be addressed during training on the wheelchair, to measure the subjects performance before and after rehabilitation, and provide a testbed and metric for the development of new technology. The program consists of three main components: the Wheelchair Skills Test (WST), a questionnaire (WST-Q), and a training program (WSTP).

In this work, we focus on the Wheelchair Skills Tests (WST) as a metric for evaluating the usability of our system for the target population. The WST consists of a set of tasks designed to test a user’s capability to operate a wheelchair. These include tasks such as moving down a hallway, parallel parking the wheelchair, moving along an incline, and rotating the wheelchair in a confined space. During evaluation, a user is asked to perform these tasks by the test administrator who determines if the task was performed successfully.

While the original test was designed for users operating power wheelchairs via a joystick, our goal is to evaluate users performing the test using the autonomous wheelchair. Use of this test provides a standardized metric for comparison between

users, as well as between use of the autonomous wheelchair and conventional control of the wheelchair.

2.6 Performance

The goal of these experiments is to determine the impact of the intelligent interface for controlling an autonomous wheelchair and its utility when used by the target population. For this thesis, we are particularly interested in the ability of the interaction manager to compensate for the noise of the speech recognition system and perform decision-making. These experiments help evaluate the POMDP as a dialogue management tool for an intelligent wheelchair. Additionally, these experiments will provide us with a set of data from a real-world domain which will be utilized in later experiments (see chapters 6 and 7). The performance of the handcrafted model used for this interface will act as a baseline for comparison against a learned model.

The performance of the system was evaluated on two test groups. Subjects were asked to execute a series of tasks from the WST using the intelligent interface to the autonomous wheelchair. Subjects were seated in the autonomous wheelchair and equipped with a microphone. When issuing a command, the subject would manually activate the microphone, say a phrase, and manually deactivate the microphone. The tests were conducted in an indoor obstacle course designed to adhere to the WST standards. For each subject, the word error rate (WER) of the recognition, number of commands issued, number of incorrect commands selected by the wheelchair, and number of queries used by the wheelchair were recorded. The transcriptions of the experiments were annotated by hand after the evaluation to obtain ground-truth

data. By examining the WER against the number of incorrect actions, the data provides a measure of the improvement of the system due to the intelligent modules.

The initial set of experiments were conducted on a set of healthy subjects, specifically, a group of occupational therapists. This group works with disabled subjects on a daily basis and has expertise in the requirements of individuals operating a powered wheelchair. This population ranged between the age of 28 to 50. Subjects were fluent in English, although many had a mild accent. The second set of test subjects are from the target population described in section 2.4. The subjects ranged in age from 31 to 85, with an average of 6.8 years of wheelchair use. While all subjects were fluent in English, their accents were stronger than the first population.

For all evaluations, the subjects were given a 20 minute training period on the wheelchair during which they were told to execute a few commands as a demonstration under the guidance of the administrator. They were then asked to perform a series of tasks unaided, as per the WST.

Tables 2-2 and 2-3 present the results of user studies on these two sets of test subject's performing the tasks of the WST. The ability of users to perform the tasks was similar between the two test groups. These results demonstrate the robustness of the intelligent interface. The subjects were able to complete the task with an average of 150 commands. In cases where the word error rate was as high as 32.7%, the system was still able to select the correct action 93% of the time. This is true across several subjects who experienced poor recognition. In these case, the system often requested additional information to help clarify the request which is preferable to issuing the incorrect action to the robot.

2.7 Discussion

This chapter presented the Smartwheeler autonomous wheelchair and the design of an intelligent interface for control. These results demonstrate the advantages of the addition of intelligent modules to process the speech input. These additional modules help reduce the error by compensating for noise in the recognition, resolving ambiguity, and reducing the complexity of the input. This is the primary motivating application for the work described in the later chapters of this thesis.

The models used in this evaluation were constructed by hand. This required extensive data collection as well as repeated trial and error. The results demonstrate that the current models used in the intelligent interface are able to handle a significant degree of error from the speech recognition system. The WER can be considered an approximation of the performance of the system which uses strictly a mapping of words to commands from the speech recognition. We can observe the impact of the decision-making component by considering the number of correct actions selected as well as the number of queries issued. Without this intelligent interface, the number of incorrectly selected actions would be closer to the WER as incorrectly recognized phrases would result in incorrect actions.

In most cases, however, the system is still selecting a few incorrect actions. While it is unrealistic to expect perfect performance, the discrepancy in the performance between users suggests improvements to the system. The performance of the system also varied per individual. Due to the effort required, creating a model specific to each individual by hand would be impractical. Finally, the models used in this evaluation are limited to the tasks defined in the WST. Expanding the model

to handle additional tasks would require repeating the model creation phase from scratch.

This suggests a method for automatic learning of the models is necessary. Learning would reduce the workload of the system developers in determining the parameters of the models, such as reward structure or the noise models of the observations used to define the interaction manager. The lengthy procedure to generate the models for this domain will have to be repeated when moving to a new domain. Ideally, the system would be able to adjust itself from experience while working with the users. This would also allow the models to adapt to each individual user, instead of using one generic model across all users. Finally, such a system would aid in the creation of new models as the domain is expanded to new tasks and domains. Developers could, for example, only be required to give rough initial models, or provide one generic model, which could adapt over time.

2.8 Related Work

The Smartwheeler autonomous wheelchair project brings together research into autonomous wheelchairs as well as from the area of dialogue management. In this section, we describe related work in these fields.

2.8.1 Autonomous Wheelchairs

Wheelchairs are a rich and promising platform for intelligent agents and are being explored by many groups as a way to help disabled and elderly individuals. Significant research has been conducted in this area, addressing the application of techniques from robots to autonomous wheelchair navigation (Simpson, 2005). The VAHM (Pruski et al., 2002) autonomous wheelchair uses a multilayer architecture

capable of navigation as well as obstacle avoidance and wall following. The selection of behaviors and modes is done manually by the user. VAHM address issues such as mapping, and includes equipping the environment with IR beacons to aid with navigation and localization. The Sharioto wheelchair (Nuttin et al., 2001) also focuses on robotic control of the wheelchair. While the original work focused on basic application of robotic navigation, recent work on Sharioto (Vanhooydonck et al., 2003) strives for a level of shared control with the user. The user directs the wheelchair via the joystick. The guidance system then attempts to infer the intention of the user using a Bayesian estimation approach to combine the environment information with the user input, and selects an action such as obstacle avoidance, passing through a doorway, or approaching a table. More recent work (Mihailidis et al., 2005) uses sensors to prevent a user from colliding with obstacles and issues voice prompts when potential dangerous situations are detected. The Hephestus (Simpson et al., 2002; Simpson et al., 2004) wheelchair project sought to develop a set of components to aid with obstacle avoidance and navigation on electric wheelchairs. The main goal was to develop these components in a way such that they could be applied to standard commercial wheelchairs with little or no modification to the chair itself. This would allow a wider adoption of such a system.

While a wheelchair can be treated as a robot, it is necessary to take the human into account. Work on motion representation with humans (Gulati & Kuipers, 2008) determines a smooth path for the wheelchair which is comfortable for the user. The results demonstrate smooth motion in areas with tight clearance that may otherwise cause the robot to slow considerable or move sharply. Work (Taha et al., 2008) has

also addressed the use of POMDPs for robot navigation and user intention. In this work, POMDPs are used to construct a model of the environment and path planning. As the user provides input via the joystick, the system attempts to determine the destination of the user, and help navigate towards it. This system also estimates parts of the model through each individual user’s experience.

Because of limitations which may arise due to subjects disabilities, research on autonomous wheelchairs has often leveraged speech recognition as an input method, though the focus of most of these projects still centers around the development of sensing and navigation.

The Vocomotion system (Amori, 1992) was early work on a speech controlled wheelchair. This system was able to recognize simple directional commands. However, this wheelchair provided no assistance with the navigation. SENARIO (Katevas et al., 1997) allowed a user to provide input via a joystick or speech recognition. In this case, the wheelchair did provide a degree of semi-autonomy to aid the user. Joystick inputs and speech were mapped to a set of commands which were executed by the wheelchair. The focus of this work, however, centered around navigation and localization. The Robchair (Pires & Nunes, 2002) autonomous wheelchair allowed the user to provide cardinal directions via the joystick or a speech recognition input. Here, a behavior-based approach was used to combine the user’s directions with obstacle avoidance and wall following behaviors. The Navchair (Levine et al., 1999) also allowed for simple speech recognition. Earlier work with the Navchair focused on aiding the user with obstacle avoidance, passage through doorways, and wall following. In later work (Demeester et al., 2006), Bayesian estimation methods

were introduced to attempt to infer the intended path of the user. By attempting to estimate the desired goal based on the command issued and sensed environment, the system could attempt to account for the uncertainty in the user’s input. Roland III (Mandel et al., 2005) focused on sensor fusion and landmark identification for mapping and navigation. Users were able to control the wheelchair through the standard joystick or the speech recognition. The speech interface on the Roland III used a semantic grammar to parse. However, the resulting policy is a “simple” mapping to actions and does not utilize a full dialogue model. The Aware Chair (Davis et al., 2003) focused on capturing the context of the user. By examining the state of the environment, such as settings of nearby devices, as well as the commands from the user, the Aware chair attempted to predict the direction of the conversation to aid with determining the user’s intention. A POMDP interface was presented for wheelchair control (Doshi et al., 2008). This work brings together dialogue management, discussed in the following section, with wheelchair control. However, this work was validated only on a small group of students.

While speech can be a powerful interface, it is not always ideal due to the task at hand or the disability of the user. Thus, research has investigated alternate means of operating wheelchairs. Wheellesley (Yanco, 2001) utilized a graphical interface on an LCD display to allow the user to select a destination. By using a vision-based navigation system, the wheelchair was able to operate in outdoor environments which may prove difficult for proximity-based sensors. A visual interface was also explored for manipulating a robotic arm attached to a wheelchair (Tsui & Yanco, 2007). By displaying the current area in front of the wheelchair on the LCD, the user is able

to select an object which is then gripped by the robotic arm. This interface was found to be faster and more comfortable than a menu based system, and reduced the complexity of object recognition by allowing the user to narrow in on the object.

EMG signals have also been used to control wheelchairs (Felzer & Freisleben, 2002). In this case, sensors are used to detect muscle contraction, which is then mapped to a set of basic commands such as “left,” “right,” “straight,” and “halt” after training with the subject. This interface is ideal for users who cannot operate a joystick or speak. Edgewise (Wobbrock et al., 2004) was developed as a gesture-based interface for input into a system. This text-input system relies on gross gestures which take advantage of the edges of the touchpad allowing even users who lack fine motor control to input text. Experiments on users found this to be a comfortable input system for text entry and designating commands to the wheelchair. Brain-computer interfaces can also be used to control a wheelchair (Philips et al., 2007). Using EEG signals, the user is able to issue the commands “forward,” “left,” and “right.” While this is a novel and potentially powerful interface, BCI control is still in early stages of development and only able to discriminate between a small set of commands.

The Smartwheeler autonomous wheelchair discussed in this thesis focuses on the development of an intelligent interface. We take advantage of existing robot navigation methods and instead focus on development and validation of a novel interface. While most voice control systems for wheelchairs focus on simple mapping of words to commands, the Smartwheeler uses full dialogue management system, including a grammar module and a probabilistic decision-making component. This

allows the interface to handle ambiguity and recognition from the speech recognition and allows the user to speak using naturally instead of learning a set of command words. Additionally, most prior work on wheelchairs validates only a small group of able individuals, often engineers or students, we worked alongside rehabilitation experts to validate the system on individuals who regularly use a wheelchair. We applied a standardized metric, in this case the WST, to allow for comparison between users and between the autonomous wheelchair and standard electric wheelchairs.

| Subject | Commands | Queries | Incorrect | % Incorrect | % Word Error |
|---------|----------|---------|-----------|-------------|--------------|
| 1 | 136 | 10 | 5 | 3.68 | 19.2 |
| 2 | 159 | 18 | 5 | 3.14 | 13.8 |
| 3 | 165 | 11 | 2 | 1.21 | 13.5 |
| 4 | 201 | 37 | 9 | 4.48 | 23.6 |
| 5 | 114 | 13 | 4 | 3.51 | 6.2 |
| 6 | 219 | 10 | 1 | 0.46 | 2.3 |
| 7 | 210 | 25 | 10 | 4.76 | 13.1 |
| 8 | 141 | 26 | 4 | 2.84 | 19.3 |

Table 2–2: Results from user studies on healthy subjects. Columns present the number of commands issued by the user, the number of queries requested by the system, the number and percentage of incorrect actions selected by the interface, and the word error rate of each subject.

| Subject | Commands | Queries | Incorrect | % Incorrect | % Word Error |
|---------|----------|---------|-----------|-------------|--------------|
| 1 | 149 | 13 | 5 | 3.36 | 12.1 |
| 2 | 145 | 20 | 6 | 4.14 | 14.5 |
| 3 | 94 | 8 | 1 | 1.06 | 7.2 |
| 4 | 122 | 24 | 2 | 1.64 | 22.9 |
| 5 | 149 | 38 | 5 | 3.36 | 30.9 |
| 6 | 117 | 2 | 0 | 0.00 | 5.2 |
| 8 | 120 | 27 | 7 | 5.83 | 23.1 |
| 10 | 149 | 35 | 10 | 6.71 | 32.7 |

Table 2–3: Results from user studies on the target population. Columns present the number of commands issued by the user, the number of queries requested by the system, the number and percentage of incorrect actions selected by the interface, and the word error rate of each subject. Subject 7 terminated the experiments early due to medication which resulted in drowsiness. The data from subject 9 contained a significant amount of noise which made hand-annotating for scoring infeasible.

2.8.2 Dialogue Management

The task of dialogue management is very similar to the interaction manager used in this project. The goal of dialogue management is to develop a robust model for human-robot dialogue which can handle the level of uncertainty in the domain. This goes beyond a simple mapping of speech recognized phrases to commands. Dialogue management faces many of the same problems as the interaction manager, as described earlier in chapter 1. Speech recognition is inherently very noisy, and the vocabulary comprises an enormous input space. The resulting dialogue should be comfortable to the user. Several approaches have been considered for dialogue management. For these approaches, the input transcriptions are considered as features. These features are used to identify a state of the dialogue which is then used to select an action.

The most basic approach to dialogue management involves the use of handcrafted dialogue managers (Balentine & Morgan, 2001; Cohen et al., 2004). Under this approach, a set of input features is mapped to a response. This mapping is created by a designer who has expert knowledge concerning the domain and must select which action to take for each possible state. This provides the advantage of a predictable, robust system. However, this requires significant effort by the designer, particularly for more complex systems. It also requires the designer to account for potential errors from the recognition system.

The handcrafted approach to dialogue management can be extended through the application of supervised learning. Here, data collected in a corpus is used to automatically determine a mapping from states to actions instead of requiring a

designer. While this reduces the workload of the designer, there are drawbacks. Dialogue has a temporal and a contextual component which is not recognized in this method. Thus, training such a system would require a specialized corpus, as examples of human to human dialogue cannot be represented (Doran et al., 2001). Instead, research has examined alternatives, such as Wizard of Oz data collection (Williams & Young, 2004) where a human operator acts as the intelligent agent providing responses from a selected set.

Markov Decision Processes (MDPs) take a similar approach to dialogue management. As with supervised learning, MDPs are able to map a set of input features to an action. However, MDPs are able to represent the sequence of the dialogue, thus preserving the context. This allows the MDP to consider the long term impact of actions and provide for more robust responses. This also allows rewards for actions and specific outcomes, so the resulting overall policy determined automatically. Reinforcement learning (RL) can also be applied to learn models and policies automatically. This type of approach will be discussed in more detail in chapter 3.1.

Research in this area of dialogue management focuses on the use of reinforcement learning to determine the correct policy. This has been seen in deployment in systems such as the NJFun phone help system (Singh et al., 2002) and the ATIS Air Travel system (Levin & Pieraccini, 1997; Levin et al., 1998). Other work has focused on reducing the required data for learning (Goddeau & Pineau, 2000), combining the use of supervised learning with reinforcement learning (Henderson et al., 2005), and considering the word error rate when selecting actions (Pietquin & Renal, 2002). MDPs, however, assume the world is fully observable and the input features provide

complete information about the world. This is particularly erroneous in light of errors in speech recognition and general ambiguity of speech.

An alternative approach attempts to consider the performance of the speech recognition system itself when selecting dialogue actions. Speech recognizers often output a confidence score representing the quality of the recognition (Kemp & Schaaf, 1997; Evermann & Woodland, 2000). This confidence score can be introduced as a new feature in addition to the transcribed words. With this additional information, a dialogue system can select which phrases to accept or reject (Pao et al., 1998; Gabsdil & Lemon, 2004) as well as request clarification. These confidence scores can also be used in coordination with cues in the dialogue, such as repeated backtracking by the user, to detect problems in the dialogue (Marc et al., 1999; Litman & Pan, 2000), and adjust the action selection accordingly.

The methods discussed so far assume a single state of the dialogue is selected based on the features, which is then used to select an action. This may leave a system vulnerable to errors in recognition. Instead, approaches have been examined which considered multiple potential states of the dialogue. Greedy decision theoretic methods (Paek & Horvitz, 2000; Paek et al., 2000) maintain a probability distribution over the potential states based on the incoming features. Actions are then selected based on the probabilities of each state and the potential reward to be received. These approaches typically focus on immediate reward, however.

Recent work has found partially observable Markov decision processes (POMDPs) to be well suited for handling dialogue management systems. As with MDPs,

POMDPs are able to maintain information about the sequence of events in the dialogue, but have the additional ability to capture the noise and uncertainty from the speech recognition and environment. This is the approach utilized for the interaction manager described in this thesis. While POMDPs are powerful models, creating the model and determining a policy has proven to be difficult. Research has focused on the development of specialized algorithms and structures to take advantage of dialogue domains and reduce the complexity of the resulting models. We will now discuss research on the use of POMDPs in dialogue management domains. POMDPs will be discussed in more detail in section 3.2.

Early work in the use of POMDPs for dialogue management presents a dialogue system for a nursing robot (Roy et al., 2000) that interacts with elderly patients in a nursing home environment. (Zhang et al., 2001a; Zhang et al., 2001b) also address dialogue management through the use of POMDPs. In this case, a Bayes net is used to process the low level input, allowing for multimodal input including speech and keyboard.

Other approaches have similarly focused on specialized structures for representing and solving POMDPs in dialogue domains. (Williams et al., 2005b) describe the integration of the speech recognition confidence scores with the POMDP resulting in a model with both discrete and continuous observations. The Summary POMDP (Williams & Young, 2005) discusses state abstraction in POMDPs by constructing a model with slots which can be filled. This approach helps limit the size of the model necessary to represent larger dialogue tasks. Later work on summary POMDPs (Williams & Young, 2006) presented a method of finding a policy on this

model. Factored POMDPs can also be used to reduce the complexity of the resulting model (Williams et al., 2005a). Under this approach, conditional independence between the state of the dialogue, the goals of the user, and the actions of the user allow simplified versions of the model to be created. A unified approach to dialogue management attempts to bring together the advantages of conventional handcrafted dialogue management systems and POMDPs (Williams, 2008). In this work, a handcrafted system selects a set of actions. The POMDP can then select one of these actions. This allows the designers to integrate domain knowledge into the system, while still allowing for the POMDP to handle for noise and ambiguity. (Doshi & Roy, 2007) consider the use of learning to adjust the model parameters of a POMDP for dialogue management, particularly for controlling an autonomous wheelchair. This work uses feedback via an oracle to help direct the learning.

The interaction manager used on the Smartwheeler is similar to the POMDP approaches to dialogue management. While dialogue management systems focus on the correct dialogue response, the interaction manager must also select from robot actions. Additionally, while dialogue management systems are concerned only with speech, the interaction manager is capable of handling multimodal input and output such as via the touchscreen display.

2.9 Future Work

The experiments described above demonstrate the first phase of the development of the Smartwheeler autonomous wheelchair. Centering the study around the Wheelchair Skills Test provides a systematic metric for evaluating the performance of the interface. The next phase would center around deployment of the wheelchair in common environments, such as a home or assisted living facility. This requires the extension of the interface to the appropriate command set based on the new tasks.

The current study relies only on speech recognition, however the intelligent interface is designed to take advantage of multimodal input. In particular, the interaction manager is able to combine information from multiple sensor sources when selecting actions. For example, the information can be presented via the touch screen display for feedback from the user. More sophisticated sensing, such as gesture recognition or eyetracking, can be brought into use through work in the long term.

Finally, the studies in this chapter only required the user to operate the wheelchair for short durations, approximately 3 hours per session. As part of deployment, long-term case studies could be conducted to determine the benefits of an autonomous wheelchair when used in everyday life. This would provide a more comprehensive understanding of the wheelchair and direct further improvements.

CHAPTER 3

Probabilistic Models for Sequential Decision-Making

Markov Decision Processes (MDPs) are powerful probabilistic models for decision-making in domains with uncertainty. According to this framework, agents move stochastically between *states* by executing *actions* and then receiving a *reward*. The agent can perform robust decision-making in light of uncertainty from the transitions. Partially Observable Markov Decision Processes extend MDPs by assuming the world is partially observable, whereas MDPs assume the state of the world is always known. Even after a stochastic action is executed, the resulting state is observed. POMDPs reduce this constraint and assume the underlying state must be inferred using *observations*. This is the case of an agent needing to use sensor information to determine its settings in the world. While POMDPs provide a much more powerful model, the increase in model complexity results in more difficult decision-making.

In this chapter we present the specifications of MDPs and POMDPs. We detail algorithms for decision-making using these models, including exact methods for determining a policy as well as approximate methods. Finally we present Reinforcement Learning (RL), a class of methods for learning and reasoning when the underlying models are not known.

3.1 Markov Decision Processes

A *Markov decision process* (MDP) is a probabilistic model of an environment used for sequential decision-making (Bellman, 1957). An agent begins in a state of the

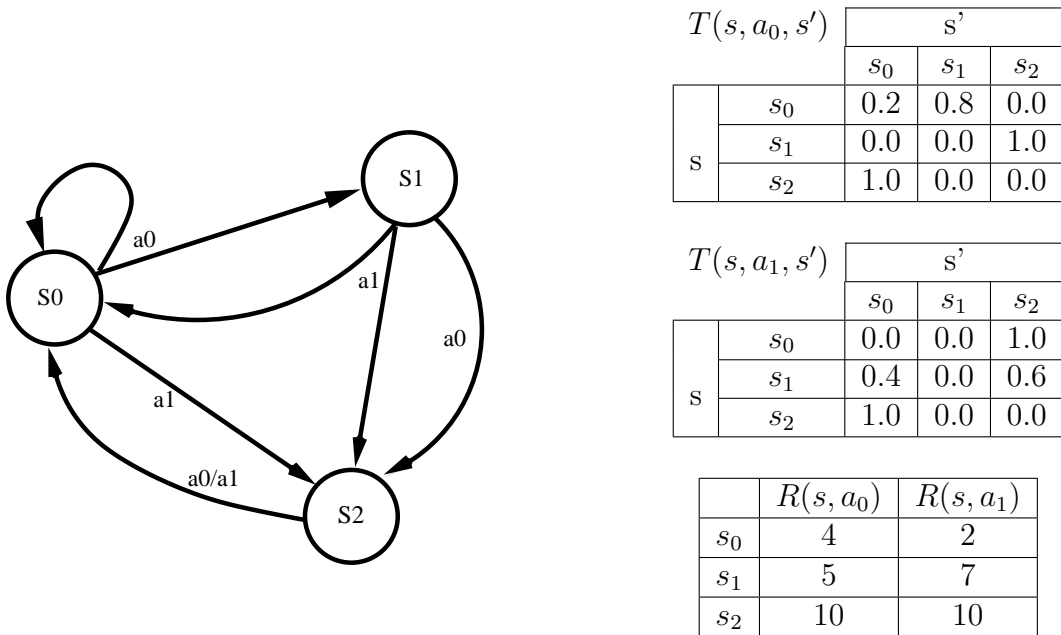


Figure 3–1: An example Markov Decision Process. In this models, s_0 , s_1 , and s_2 are states. a_0 and a_1 are actions.

world, s , and executes an action, a . This action leads the agent non-deterministically to a new state, s' , at which point the agent selects a new action. An example MDP can be seen in figure 3.1

Formally, an MDP is defined as a tuple (S, A, R, T) where

- $S = \{s_0, s_1, s_2, \dots, s_n\}$ is a set of states of the world
- $A = \{a_0, a_1, a_2, \dots, a_m\}$ is a set of actions executable by the agent
- R is the reward function which determines the reward an agent receives when executing an action, a , in a state, s .

$$R(s, a) : S \times A \rightarrow \mathbb{R} \tag{3.1}$$

- T is the transition function defining how an agent moves between states based on the properties of the environment. One of the primary principles behind MDPs is the *Markov property*. A process is said to be *Markovian* if the future states are conditionally dependant only on the current state. This assumption allows the transitions to be defined only over the current state and action selected:

$$T(s, a, s') := Pr(s_t = s' | s_{t-1} = s, a_{t-1} = a) \quad (3.2)$$

Given a model of the environment, a policy can be found to allow the agent to select actions. This *policy*, π , provides a mapping from states to actions.

$$\pi : s \rightarrow a \quad (3.3)$$

An *optimal policy* is one which maximizes the value of each state, $V(s)$, based on the *expected cumulative reward* of the agent. This is know as the *value function*:

$$V^*(s) = E_\pi \left[\sum_{t=0}^T \gamma^t r_t \right] \quad (3.4)$$

where r_t is the reward received at time t , $\gamma \in [0, 1)$ is the discount factor, and $T \in [1, \infty)$ is the planning horizon. The discount factor reflects the idea that immediate rewards are more significant than future rewards. The smaller the value of the discount factor, the less influence future actions will have on the current decision. This implies that the value of a state is based on the relationship to the successor states. This value function can be defined recursively as follows:

$$V^t(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^{t-1}(s') \right] \quad (3.5)$$

This is known as the Bellman equation and acts as the basis for many MDP solution techniques.

A *fixed point* occurs in the iteration when the value function no longer changes between iterations. This defines the optimal value function, V^* :

$$V^*(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s') \right] \quad (3.6)$$

3.1.1 Value Iteration

Value Iteration is a dynamic programming approach to solving MDPs based on the Bellman equation formulation. By using a recursive approach, value iteration repeatedly updates the value of each state for $t = 1, 2, \dots$ until a terminating condition is met, usually when there is little or no change in the value function. Value iteration is presented in Algorithm 1.

Algorithm 1 Value Iteration

```

 $V_1(s) = 0 \quad \forall s \in S$ 
 $t = 1$ 
repeat
   $t = t + 1$ 
  for all  $s \in S, a \in A$  do
     $V_t(s) = \max_{a \in A} [R(s, a) + \gamma [\sum_{s' \in S} T(s, a, s') V_{t-1}(s')]]$ 
  end for
until  $|V_t(s) - V_{t-1}(s)| < \epsilon \quad \forall s \in S$ 

```

The terminating condition is controlled by parameter ϵ . The fixed point occurs when ϵ is 0. However, iterating until this condition is met is often impractical. Instead, the fixed point can be approximated by setting ϵ to small values. The closer ϵ is to 0, the better the approximation. However this will result in more iterations and

computation. Typically this value is determined based on the domain and resources available.

MDPs have proven to be powerful tools, but are of limited use in robotics applications due to their assumption that the world is fully observable. Other research centers on topics such as determining policies for very large MDPs (Meuleau et al., 1998), modelling and decision-making in continuous domains (Szepesvári, 2001), and the use of MDPs in multiagent systems (Guestrin et al., 2001). In this thesis, we focus on partially observable domains. Further information concerning MDPs can be found here (Puterman, 1994; Sutton & Barto, 1998) but is beyond the scope of this work.

3.2 Partially Observable Markov Decision Processes

One of the key assumptions of MDPs is that the environment is fully observable. While the agent may not know the outcome of executing an action beforehand, the agent is aware of the resulting state after execution. While this may be a valid assumption in some domains, many agents only have limited or incomplete information about the environment. This is known as a *partially observable* domain. Instead, these agents must infer the state based on *observations* of the environment. *Partially observable Markov decision processes* (POMDPs) (Sondik, 1971; White, 1991; Lovejoy, 1991; Kaelbling et al., 1998) are an extension to MDPs which account for the noise and ambiguity due to sensing the environment. As with MDPs, an agent begins in a state, which may not be known, and executes an action. The agent stochastically transitions to a new state, and receives an observation. The new state

is not known, instead, the agent uses the transition and observation model to update its belief about the world, and selects the next action.

Formally, a POMDP is defined as a tuple (S, A, Z, R, T, O, b_0) where

- $S = \{s_0, s_1, s_2, \dots, s_n\}$ is a finite set of states of the world.
- $A = \{a_0, a_1, a_2, \dots, a_m\}$ is a finite set of actions executable by the agent.
- $Z = \{z_0, z_1, z_2, \dots, z_q\}$ is a finite set of observations observable by the agent.
- R is the reward function which determines the reward an agent receives when executing an action, a , in a state, s :

$$R(s, a) : S \times A \rightarrow \mathbb{R} \quad (3.7)$$

- T is the transition function defining how an agent moves between states based on the properties of the environment:

$$T(s, a, s') := Pr(s_t = s' | s_{t-1} = s, a_{t-1} = a) \quad (3.8)$$

- O is the observation function defining how agent receives observations from the environment:

$$O(s, a, z) = Pr(z_t = z | s_t = s, a_t = a) \quad (3.9)$$

- b_0 is the initial distribution over the states.

3.2.1 Belief State Tracking

Because the system is partially observable, the agent does not generally observe the true state of the world. Instead, a *history* maintains a trace of the agent's actions and observations over time. The history at time t is defined as:

$$h_t = a_0, z_1, a_1, z_2, \dots, a_{t-1}, z_t \quad (3.10)$$

Maintaining an extensive history quickly becomes impractical. Instead, a distribution, known as a *belief state*, is maintained over the set of states, S . Let $b_t(s)$ denote the probability of an agent being in state s at time t . The set of all possible beliefs is defined as ΔS . This simplex is a continuous space of dimension $|S| - 1$.

Given an action, a_t , and resulting observation, z_{t+1} , the belief can be updated as follows:

$$b_{t+1}(s') = \frac{O(s', a_t, z_{t+1}) \sum_{s \in S} T(s, a_t, s') b_t(s)}{Pr(z_{t+1} | a_t, b_t)} \quad (3.11)$$

where

$$Pr(z_{t+1} | a_t, b_t) = \sum_{s \in S} \sum_{s' \in S} O(s', a_t, z_{t+1}) T(s, a_t, s') b_t(s) \quad (3.12)$$

is the probability of observing z after execution action a_t in belief state b_t . This acts as a normalizing factor such that the probabilities over the beliefs sum to 1. Maintaining the belief state allows the agent to compress all the information about the history into a compact sufficient representation. At each time, $t + 1$, as the agent executes an action, a , and receives an observation, z , the belief state is updated to reflect this information based on the belief state at time t . Furthermore, it has been shown that the belief state is a sufficient representation for optimal decision-making in POMDPs (Sondik, 1971).

3.2.2 Exact Value Iteration

The goal of decision-making in a POMDP is to determine an optimal policy, π^* , which acts as a mapping from belief states to actions.

$$\pi^* : b \rightarrow a \quad (3.13)$$

As an agent executes actions and the belief state is updated, the POMDP policy dictates the next action to take for the updated belief state. This policy should maximize the value function, as defined as the sum of discounted future rewards:

$$V^\pi(b) = E_\pi \left[\sum_{t=0}^T \gamma^t r_t | b_0 \right] \quad (3.14)$$

where r_t is the reward received at timestep t , γ is the discount factor, and $V^*(b)$ is the value of the executing the optimal policy starting at belief state, b_0 . The value of each belief state is then defined via expectation over the underlying states:

$$V^*(b) = \sum_{s \in S} b(s) V^*(s) \quad (3.15)$$

A form of value iteration can be used to find the optimal policy. The value function for planning horizon $t + 1$ can be determined iteratively as follows:

$$V_0(b) = \max_{a \in A} \sum_{s \in S} b_0(s) R(s, a) \quad (3.16)$$

$$V_{t+1}(b) = \max_{a \in A} \left[\sum_{s \in S} b(s) R(s, a) + \gamma \sum_{z \in Z} Pr(z|a, b) V_t(b_{a,z}) \right] \quad (3.17)$$

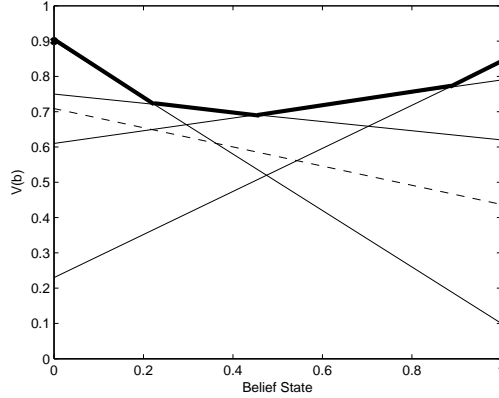


Figure 3-2: Example value function. The dark lines show the optimal value. Each vector is associated with an action. The dashed line is dominated by other vectors and has no influence on the resulting policy.

where $b_{a,z}$ is the resulting belief state after executing a and observing z , and $R(b, a)$ is the expected reward of executing a in belief state b , as defined by:

$$R(b, a) = \sum_{s \in S} b(s) R(s, a) \quad (3.18)$$

The optimal policy at $t + 1$ can then be expressed based on the value function:

$$\pi_{t+1}^*(b) = \arg \max_{a \in A} \left[\sum_{s \in S} b(s) R(s, a) + \gamma \sum_{z \in Z} Pr(z|a, b) V_t(b_{a,z}) \right] \quad (3.19)$$

It has been shown that the optimal value function is a piecewise linear, convex function over the continuous belief state and can be represented as a set of vectors. The Enumeration algorithm (Sondik, 1971) was presented as a method for determining this value function. For a given planning horizon t , the value function can be represented as a set of α -vectors:

$$\Gamma_t = \{\alpha_0, \alpha_1, \dots, \alpha_m\} \quad (3.20)$$

Each of these α -vectors is associated with one action and represents the expected reward of executing that action and continuing to follow the policy. The value function can now be expressed in terms of these α -vectors:

$$V_t(b) = \max_{\alpha \in \Gamma_t} \sum_{s \in \mathcal{S}} \alpha(s)b(s) \quad (3.21)$$

The optimal policy can be determined by selecting the action associated with the α -vector which maximizes the value function at b .

An α -vector is said to be dominated if:

$$\forall b \in \Delta S \exists \alpha' \in \Gamma \text{ s.t. } \sum_{s \in S} \alpha'(s)b(s) \geq \sum_{s \in S} \alpha(s)b(s) \quad (3.22)$$

Such an α -vector will never be part of the optimal policy and can be removed from the set Γ_t . Similarly, α_1 is said to dominate α_2 if and only if $\alpha_1(s) \geq \alpha_2(s) \forall s$.

An iterative approach can be used to calculate the set Γ_t . We expand equation 3.17 and express the value function in terms of the α -vectors:

$$V_{t+1}(b) = \max_{a \in A} \left[\sum_{s \in S} b(s)R(s, a) + \gamma \sum_{z \in Z} \max_{\alpha \in \Gamma_t} \sum_{s \in S} \sum_{s' \in S} T(s, a, s'), O(s, a, z)V_t(b_{a,z}) \right] \quad (3.23)$$

We initialize the set of α -vectors for a single step horizon, $t = 1$, as follows:

$$\Gamma_1^a = \{\alpha^{a,*} | \alpha^{a,*}(s) = R(s, a)\} \quad (3.24)$$

The α -vectors for $t + 1$ can now be constructed recursively using the set from t .

We first construct intermediate sets $\Gamma_{t+1}^{a,*}$ and $\Gamma_{t+1}^{a,z}$:

$$\Gamma_{t+1}^{a,*} = \{\alpha^{a,*} | \alpha^{a,*}(s) = R(s, a)\} \quad (3.25)$$

$$\Gamma_{t+1}^{a,z} = \left\{ \alpha^{a,z} | \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s, a, s') O(s', a, z) \alpha'(s') \quad \forall \alpha' \in \Gamma_t \right\} \quad (3.26)$$

Essentially, we are considering all the possible outcomes of observations, z , when executing action, a , and selecting α -vectors from the previous iteration. Then a second intermediate set is generated by combining the first intermediate sets using a cross-sum¹ operator:

$$\Gamma_{t+1}^a = \{\Gamma_{t+1}^{a,*} \oplus \Gamma_{t+1}^{a,z_1} \oplus \Gamma_{t+1}^{a,z_2} \oplus \dots \oplus \Gamma_{t+1}^{a,z_{|Z|}}\} \quad \forall a \in A \quad (3.27)$$

We construct the final set of vectors for $t + 1$:

$$\Gamma_{t+1} = \bigcup_{a \in A} \Gamma_{t+1}^a \quad (3.28)$$

This algorithm can be used to find an exact policy for a POMDP. Unfortunately, the computation time for each iteration grows exponentially (Kaelbling et al., 1998). Let $|A|$ be the number of actions, $|S|$ be the number of states, $|Z|$ be the number of

¹ The cross-sum (\oplus) operator is defined: for sets $A = \{a_1, a_2, \dots, a_{|A|}\}$, $B = \{b_1, b_2, \dots, b_{|B|}\}$, $A \oplus B = \{a_1 + b_1, a_1 + b_2, \dots, a_1 + b_{|B|}, a_2 + b_1, \dots, a_{|A|} + b_{|B|}\}$

observations, and $|\Gamma_{t-1}|$ be the number of α -vectors in the previous iteration. Equation 3.25 generates $|A||Z||\Gamma_{t-1}|$ vectors. The cross-sum in equation 3.27 generates $|A||\Gamma_{t-1}|^{|Z|}$ vectors. Thus, the set of α -vectors for horizon t can be computed in time $O(|S|^2|A||\Gamma_{t-1}|^{|Z|})$. This makes exact planning via the Enumeration algorithm for even moderately sized POMDPs unreasonable. Therefore, research has been conducted into methods to approximate the policy allowing policies for larger models to be found.

Several approaches have been proposed to determine exact solutions for POMDPs. The Enumeration algorithm is considered to be the most straightforward. Subsequent research attempts to improve on the performance.

The One-pass algorithm (Sondik, 1971) attempts to construct the set of α -vectors by identifying regions where the vectors are dominant. This algorithm selects belief points at random, calculates the associated α -vectors, and attempts to determine regions where the vector is guaranteed to be dominant. The linear support algorithm (Cheng, 1988) also attempts to identify regions of the belief where a vector is dominant. To accomplish this, the linear support algorithm focuses on points at the intersection of two vectors based on the idea that the error in the value function will be greatest at the edges. This algorithm begins by selecting a point and generating the α -vector for the point. The points at the edges of the vector, where it intersects other vectors, are examined and the α -vectors calculated. If the new α -vectors return a higher expected reward, the process is repeated.

The Witness algorithm (Littman, 1996) takes a similar approach. As with the linear support algorithm, the Witness algorithm first determines the vector for a

belief point, and then attempts to determine where the vector is not dominant. Witness considers different observations to determine if vectors with higher values exist.

Incremental pruning improves on the Enumeration algorithm by attempting to reduce the number of vectors at each iteration. During the final step of the vector construction, in equation 3.27, incremental pruning eliminates vectors which are dominated by other vectors. This prunes out vectors that will never be selected in future iterations (Cassandra et al., 1997).

3.2.3 Point-based Approximate Methods

Because exact decision-making for POMDPs is largely intractable, significant research has been conducted on approximate methods. They allow a tradeoff between quality of the policy and availability of solution.

Sampling of the belief state has become a popular method for determining policies for POMDPs (Lovejoy, 1991; Hauskrccht, 2000; Pineau et al., 2003). These methods typically function by determining the value function over only the belief states that are reachable. This is accomplished by only performing value iteration over a set of samples of the belief states. They have the advantage that by only operating on a set of belief states at each iteration, resource use is limited. As with the exact methods, they typically behave as anytime algorithms, allowing a solution to be returned at any point during execution. The size of the sample set also handles some of the memory constraints problems that may be encountered using exact methods.

Point-based approximation methods operate over a set of belief points, B , where B is a subset of the belief simplex ΔS . This subset represents a set of reachable belief

states from the initial belief state, b_0 , of the model. To find a policy over a POMDP using these methods, a *backup* operator is introduced:

$$\text{backup}(b) = \arg \max_{\alpha \in \Gamma_t} \sum_{s \in S} b(s) \alpha(s) \quad (3.29)$$

The *backup* operator essentially selects the best α -vector for a belief state given a set of potential vectors Γ_t . As with the exact methods, we first initialize the set Γ_1^a :

$$\Gamma_1^a = \{\alpha^{a,*} | \alpha^{a,*}(s) = R(s, a)\} \quad (3.30)$$

For each iterative step $t + 1$, we then calculate sets $\Gamma_{t+1}^{a,z}$ as before:

$$\Gamma_{t+1}^{a,z} = \{\alpha^{a,z} | \alpha^{a,z}(s) = \gamma \sum_{s' \in S} T(s, a, s') O(s', a, z) \alpha'(s') \quad \forall \alpha' \in \Gamma_t\} \quad (3.31)$$

and a second intermediate set $\Gamma_{t+1,b}^a$:

$$\Gamma_{t+1,b}^a = \left\{ \Gamma_1^a + \sum_{z \in Z} \arg \max_{\alpha \in \Gamma_{t+1}^{a,z}} \left[\sum_{s \in S} \alpha(s) b(s) \right] \right\} \quad (3.32)$$

This effectively marginalizes over the observations. As the final step in the iteration, we apply the *backup* operator as before:

$$\text{backup}(b) = \arg \max_{\alpha \in \Gamma_{t,b}^a} \sum_{s \in S} b(s) \alpha(s) \quad (3.33)$$

This defines a general approach for a point-based solver for a POMDP. Research has presented different algorithms which take different approaches to some of the aspects of the generic algorithm.

The point-based value iteration (PBVI) (Pineau et al., 2003) algorithm is one example of this solution methodology. At each iteration, the value function is determined at a selected set of points of the belief space. After each iteration, new belief points are selected by simulating each action followed by resampling a random observation. The L_1 norm distance is calculated between the new belief states and belief states in B , discarding new beliefs which already exist in the sampled set, and adding the new beliefs furthest from the set. At each iteration, one belief is added for each existing belief, at most doubling the size of the B . A similar approach is seen in the Perseus algorithm (Spaan & Vlassis, 2005). Here, a large set of samples of the belief state are select intially and maintained. Only a subset of these samples is updated at each iteration, the idea being that one update can affect many samples. At each iteration, Perseus updates the α -vector for a random belief from B , and notes which other beliefs have also been improved by the new α -vector. This is repeated until all beliefs in B have been improved. This reduces computation as α -vectors from one belief state may also yield an improved value for other belief states, thus not requiring a calculation at those states.

Heuristic Search Value Iteration (HSVI) (Smith & Simmons, 2004) focuses on the selection of belief points. This algorithm uses a heuristic search through reachable belief states to select the set of belief points. HSVI maintains estimates of the upper and lower bounds of the value function and attempts to select belief points to reduce the difference of the bounds. These bounds are updated when a backup operation occurs. Forward Search Value Iteration (FSVI) (Shani et al., 2007) extends the work of HSVI to improve the belief state selection heuristic. FSVI assumes the POMDP is

fully observable, determines the policy of the resulting MDP, and uses that policy to guide a search through the belief space to select points. Successive Approximation of the Reachable Space under Optimal Policies (SARSOP) (Kurniawati et al., 2008) also takes a similar approach as HSVI. SARSOP attempts to sample the belief points which would be reached by executing the optimal policy. Because the optimal policy is not known, the bounds on the value function are used to approximate the policy.

Details of these planning procedures are beyond the scope of this thesis. One step in the framework to be discussed in chapter 5 requires solving POMDPs. The choice of method is not critical to the framework. The work presented in this thesis utilizes SARSOP as it was found to be the fastest algorithm at the time, although earlier versions of our work have used PBVI and Perseus. In practice, any exact or approximate POMDP solution method could be used.

3.3 Reinforcement Learning

Up to this point, this chapter has focused on planning methods for MDPs and POMDPs. However, all of the techniques discussed assume the parameters of the model of the world are provided. For cases where the model is not known, a class of learning algorithms known as *Reinforcement Learning (RL)* (Sutton & Barto, 1998) exists. Reinforcement learning allows an agent to learn about the environment while executing actions. In this section, we present some of the fundamental research on reinforcement learning approaches. We describe some initial work in the field focused on fully observable MDPs, and continue with work in partially observable environments. We distinguish between methods which learn only the policy and those which also learn the underlying models.

3.3.1 Reinforcement Learning in Fully Observable Domains

Temporal difference learning, $TD(\lambda)$ (Barto et al., 1983) is an early example of a model-free reinforcement learning method. $TD(\lambda)$ attempts to update the value function of each state, $V(s)$, as the agent executes a policy through the environment and observes the resulting state and reward. The policy can then be updated using the estimates. Q-learning (Watkins, 1989) extends the work in reinforcement learning by providing a unified approach. Q-learning estimates the expected reward for executing an action in a state, $Q(s, a)$ based on the reward received and value of the resulting state. This allows the policy to be directly learned by selecting the action with the maximum Q value in each state. A similar algorithm known as State Action Reward State Action (SARSA) (Rummery & Niranjan, 1994). SARSA differs from Q-learning by considering the action taken by the agent in the following steps, while Q-learning considers the action which is expected to be taken based on the policy.

$TD(\lambda)$ and Q-learning are both *model-free* methods as they only learn a mapping from states to value functions. However, there are cases where learning a model of the world is desirable. A simple approach to *model-based* learning would be to explore the environment, learn the transition and reward parameters, and then determine a policy using this model via value iteration. This straightforward approach poses several problems. First, the choice of policy to gather the data is not defined. An undirected exploration of the environment would require significantly more data than one which was able to focus the search on relevant parts of the space. Second, this approach would fail to account for dynamic environments. Any changes to the environment after the learning phase is complete would be unaccounted for and

problematic. To address these issues, model-based reinforcement learning methods have been developed to interleave the learning and decision-making phases during execution and allow for continuous learning.

Dyna (Sutton, 1990) introduces an algorithm for model-based reinforcement learning. At each step, after transition from s to s' via action a , Dyna first updates the model by maintaining statistics about the transitions and rewards received. Because the environment is fully observable, simply maintaining counts is sufficient to estimate these parameters. Using the newly estimated parameters, the policy can then be updated with an iteration of the Bellman equations. This update occurs to a set of randomly selected states in addition to s . One of the weaknesses of Dyna centers around the random selection of points to update each step. Prioritized Sweeping (Moore & Atkeson, 1993) was introduced to help direct the update phase of the algorithm. Prioritized sweeping maintains a priority for each state and uses this priority to order the updates. This is accomplished by maintaining the estimated value of each state, $V(s)$, and considering the change in the estimate after an update as a metric for the priority on the state's predecessors. Conceptually, when one state is updated, it implies that the states which lead to that state also need to be updated.

There exists extensive literature exploring the use of reinforcement learning in fully observable domains (Kaelbling et al., 1996; Sutton & Barto, 1998; Szepesvari, 2010). In this thesis, we instead focus on learning in partially observable domains.

3.3.2 Learning in Partially Observable Domains

One approach to learning in partially observable domains is to operate directly on the observed information. Nearest Sequence Memory (NSM) (McCallum, 1994)

uses this approach to perform reinforcement learning. As with Q -learning, NSM attempts to learn Q -values. However, these values are learned over instances instead of underlying states. Using NSM, when an agent takes a step in the environment, it records the action taken, and observation and reward received as an instance. Using a similarity metric, which can vary, NSM then identifies the k nearest neighbors from previously encountered instances. The value for each potential action at the current instance is calculated by average the value of that action in the neighbors and selecting the action with the highest value. The Q -values for the nearest neighbors are then updated based on the outcome of the action.

While NSW works with an instance, Utile Suffix Memory (USM) (McCallum, 1995) takes a similar approach but using a history instead of a single instance. USM attempts to cluster an instance based on its history against other instances with similar histories. To perform this clustering, USM first constructs a suffix tree during a training phase where data is collected. The root node of the tree contains all the possible instances. The children of the root each represent one observed observation, with each node containing the instances that contain the observation at $t - 1$, as well as the Q -values for each action from that instance. The next layer of the tree represents the actions which were executed at $t - 1$. Again, the nodes each contain the instances which executed that action, then observed the observation of the parent node. The tree continues to alternate between action and observation layers, with each layer further dividing the instances based on the history. This results in instances in common nodes at deeper layers sharing a longer history. When considering a new instance, the history of the instance is traced through the tree beginning

at the root until a leaf node is reached. If the Q -values are statistically different, the node is split. Otherwise, the new instances are added to the existing node. For action selection, when a new instance is encountered, again the history is traced through the tree until a leaf node is found. At this point, the Q -values of the instances in the node are averaged and used to select an action. Once an action is executed, the values of the instances are updated.

Noisy Utile Suffix Memory (NUSM) (Shani & Brafman, 2004) extends the initial work on USM to take advantage of situations where the observation model is known. In this work, it is assumed $O(z|s)$ is known in advance. To integrate the observation information, an additional layer is first added to the tree representing the potential states. When a new instance is considered, a weight is calculated based on the observations and the trace is inserted into multiple nodes in the tree. This allows the system to consider the probability of a observation sequence occurring when performing the clustering. Later work (Shani et al., 2005) demonstrates the ability to convert the tree from the NUSM algorithm into a POMDP. To provide for this, the leaves of the tree are the resulting states, the rewards are obtained from the Q -values of the instances in the nodes, and the transitions are determined by maintaining counts during the training phase. However, this work was unable to learn the observations and assumes they are determined offline using sensor models.

The idea of utilizing the observed information for reinforcement learning has also been addressed through the application of SARSA (Loch & Singh, 1998). Here, at each time step, the tuple of observed information, (s_t, a_t, r_t, s_{t+1}) is treated as a

state. SARSA is then executed over this new representation to learn a policy for the environment as if it were a fully observable domain.

Research has also shown that a stochastic policy over the observations can be applied to perform learning and decision-making in POMDPs (Jaakkola et al., 1995). While other research determines policies over instances and histories, which include actions executed, observations, and rewards received, this work determines a mapping from observations to actions through a stochastic policy. To evaluate the policy, a Monte Carlo sampling approach is applied to generate traces through the environment and estimate the value of each observation-action pair, $Q(z, a)$. Using these updated estimates, the stochastic policy is updated and the process repeated.

One approach to learning in POMDPs presents a *homing policy* to aid with the learning (Even-Dar et al., 2005). In this approach, a policy is found which returns the agent to a fixed belief state, essentially “resetting” the POMDP. This allows the agent to execute sequences of actions from the same starting point. Using this homing policy, algorithms were presented to learn the policy directly as well as learn the transition and observation parameters while utilizing this homing policy.

Recent work has proposed an alternative representation for partially observable domains known as Predictive State Representation (PSRs) (Littman et al., 2001; Singh et al., 2004). PSRs attempt to represent the environment using a series of *tests* based on the actions taken and observations received, eliminating the need to model the underlying POMDP of the environment. While it has been demonstrated that PSRs are equivalent to POMDPs for representation (Littman et al., 2001), they

require a new set of methods for planning (Izadi & Precup, 2003) and learning (Singh et al., 2003).

Recent work on reinforcement learning in partially observable domains has focused around Bayesian techniques. These approaches are discussed in details in chapter 4.

3.4 Summary

In this chapter, we presented a set of probabilistic models for decision-making in stochastic environments. Markov Decision Processes are able to model noise from transition. However, they assume the environment is fully observable. Partially Observable Decision Processes relax this constraint by introducing observations to model noise from sensing or aliasing in general. We presented exact algorithms for decision-making using these models. In the case of POMDPs, approximate methods are also needed to allow solutions to larger models to be found. However, all of these techniques assume models of the environment exist. To handle situations where models are unknown, reinforcement learning algorithms present approaches for determining policies and models from experience. Reinforcement learning techniques assume the world is fully observable and research into partially observable domains is limited. Learning in partially observable environments is crucial as agents are used in more complex domains, such as the interaction manager described in chapter 2, where inputting models by hand becomes less feasible.

CHAPTER 4

Bayesian Reinforcement Learning

In the previous chapter, we defined various Markov models as well as detailing techniques for decision-making. These models provide a rich framework for many real world domains due to their capability to capture and represent uncertainty. However, performing decision-making assumes that a model of the environment is provided.

Specifying the model of a problem can be a difficult task, especially as problems become extensive and more complex as often is the case in real-world environments. In domains with many states, the designer may not have complete information about the entire environment, even in cases where redundancy allows states to be generalized. This is also true of agents with many actions and observations. Proper models of the system dynamics may be missing. Even when available, designers are required to have extensive knowledge of the domain to craft acceptable models, which is often time consuming. This domain knowledge may involve specialized techniques to model components, such as sensor models for laser scanners, error rates for recognition systems, and physics models for robot actuators. These techniques are often custom to the sensor. A system capable of automatically determining these models would be robust to the addition and removal of these types of components. A similar problem arises due to dynamic environments which may change over time. This may be a result of a continuously changing environment, or a situation where

models are created for one domain and later applied to a separate, though similar, domain. Finally, many problems may not be easily expressed explicitly by the designer. A system able to determine its own model parameters may be able to find a more appropriate representation which may be less intuitive and human readable, but more adept for recognition and decision-making.

In section 3.3, several techniques were presented for estimating model parameters and performing decision-making based on empirical data. However, there are still a variety of concerns, especially when moving to a partially observable domain. Reinforcement learning algorithms such as Q-learning and Sarsa assume a fully observable environment. Instance-based methods, such as NSM and USM are able to learn policies for partially observable domains. However, they assume the reward is an observed parameter which is often not realistic in real world domains. The issue of exploration and exploitation also remains. Additionally, these are model-free methods and many situations exist where a model would be desirable. Finally, expectation-maximization algorithms are ideal for learning transition and observation parameters, but require extensive training data and have no methods for learning the reward parameter.

In this chapter, we review a class of approaches to learning known as Bayesian Reinforcement Learning (BRL). Bayesian reinforcement learning is an extension of the reinforcement learning techniques described in chapter 3. However, instead of using point estimates of the parameters, Bayesian reinforcement learning represents and maintains distributions over the unknown parameters. We present an overview

and general formulation of BRL. We then review the literature concerning the application of BRL to Markov Decision Processes and the extensions into partially observability.

4.1 Overview

In probability theory, *Bayes Rule* defines the relationship between a conditional probability over two events. This relationship is defined as follows:

$$P(A = a|B = b) = \frac{P(B = b|A = a)P(A = a)}{P(B = b)} \quad (4.1)$$

where A and B are random variables which take on values a and b respectively. The goal is to determine the probability that $A = a$ based on a second observed value, $B = b$. This *posterior probability*, $P(A = a|B = b)$, can be updated based on the *likelihood*, $P(B = b|A = a)$, and the priors over the individual variables, $P(A = a)$ and $P(B = b)$. When applied, Bayes Rule allows the belief that an event occurred to be updated based on new evidence. If $B = b$ was an observed event, the belief that $A = a$ can then be updated.

Algorithm 2 A general framework describing Bayesian reinforcement learning. Various approaches address the challenges in each step through different techniques.

Initialize distributions over unknown parameters

loop

 Select action based on distributions

 Execute action

 Observe resulting state, reward, and/or observation

 Update counts of unknown parameters based on observed information

end loop

The aim of Bayesian reinforcement learning (BRL) is to maintain a posterior distribution over possible model parameters and to compute an action selection policy which is optimal with respect to this posterior. A general formulation for Bayesian RL is shown in Algorithm 4.1. First, distributions are initialized over the unknown parameters in the models. Using the current information about those parameters, an action is selected. After the action is executed, the agent takes the resulting information about the environment and updates the distributions over the parameters.

This general formulation of BRL is defined by three main issues: representation of the data, updating the parameters given data, and action selection. Representation and updating are common to any Bayesian learning approach. Action selection is specific to BRL.

Representation: One of the fundamental questions when using Bayesian reasoning is the choice of representation. A common approach is to represent the unknown parameters with a set of distributions defined by hyper-parameters. This provides a compact representation which captures information about the parameter values themselves as well as the uncertainty over the parameters. The choice of distribution can also reduce the complexity by allowing a closed-form update of the distribution. As parameters such as transitions and observations are multinomials, a natural choice of distribution is a Beta or Dirichlet distribution, which act as conjugate priors to the multinomials. Distributions such as these, where the hyper-parameters represent counts of a events, allow for a more simple update as the counts can be incremented when evidence is observed. Finally, methods in BRL usually take advantage of the conjugate prior of distributions. For example, as the transition distribution for an

MDP is a multinomial, it is a natural fit to represent the prior distribution as a Dirichlet. The Beta distribution is a second common representation as the conjugate prior of a binomial distribution.

Parameter Update: Given a representation, the system must be able to update the hyper-parameters based on experience. In the case of fully observable environments, data which is observed often has an associated hyper-parameter which can be incremented directly. These hyper-parameters are related to the states. Alternate approaches must be considered for partially observable environments, as the underlying state is not witnessed directly, instead relying on observations.

Action Selection: The algorithms described in chapter 3 for determining policies over MDPs and POMDPs assume point estimate representations of the parameters. When the parameters are represented using distributions as opposed to point estimates, algorithms such as value iteration no longer become directly applicable. Instead, methods for action selection are needed which reason over the distributions, taking into account the uncertainty in the current parameters. Work in BRL has presented exact approaches which can be used for selecting optimal actions based on the distributions (See sections 4.2 and 4.3). Because exact solutions can be costly to calculate, often sampling methods are used to approximate the optimal policy by drawing representative examples from the distributions. Once these models are sampled, an algorithm for determining the policy as discussed earlier, such as value iteration, can be used to determine a policy.

Bayesian reinforcement learning maintains many of the advantages found in traditional reinforcement learning algorithms. An agent is able to execute actions and

operate in the environment while concurrently learning a model or policy. This interleaving of learning and execution removes the necessity to gather extensive training data as a separate procedure and allows the agent to progress towards task completion immediately after initialization.

This type of approach also allows the agent to select actions based on the current state of the learning process. The agent can use information, such as the uncertainty over the learned parameters, to help select an optimal action. In a similar fashion, the agent can select actions to explore parts of the search space it deems useful, essentially selecting actions specifically to gain information and learn parts of the model which may be lacking data.

Finally, by allowing learning to occur during execution, the agent is able to adapt to dynamic environments. As the environment changes, reinforcement learning allows the agent to adjust the model as necessary, as the new experiences are encountered. In the case of the dialogue manager, a user's preferences of the system behavior may change over time. By providing a continuous mechanism for learning, the system can continue to adjust the models based on the preferences of the user.

4.2 Bayesian Reinforcement Learning in Fully Observable Domains

One of the earliest examples of Bayesian reinforcement learning demonstrated the application of Bayesian reasoning in a model-free Q-learning system (Dearden et al., 1998). In this work, distributions are used to maintain the uncertainty in the Q values for an agent executing Q-learning. By considering this uncertainty when performing action selection, the agent is able to direct the search towards poorly known regions of the learning space. This helped address the exploration

versus exploitation problem in learning during execution. Different approaches are presented to select the actions based on the current distributions. These approaches relied on heuristics such as expected reward of each current action, the information to be gained, as well as a technique which samples a set of actions directly from the distributions.

Similar work (Dearden et al., 1999) first extended Bayesian reinforcement learning to model-based methods using Markov decision processes and provides many of the basic ideas used in model-based learning. In this work, the goal was to maintain a distribution over the potential models themselves as a mechanism for determining the Q-values and selecting actions. This is accomplished by maintaining a Dirichlet distribution over the unknown parameters, in this case, the transition and reward models. These Dirichlet distributions were defined by a set of hyper-parameters, $N(s, a, s')$ and $N(s, a, r)$. These represented counts over the number of times a transition from s to s' occurred after executing action a and the number of times a reward r was received after executing action a in state s . As an agent experienced a transition or reward, the appropriate count is incremented. Based on the set of experiences, μ , these Dirichlet distributions can be used to define the posterior distribution over the potential models, $P(M|\mu)$. Using this, a belief over the potential models can be generated.

To perform action selection based on the posterior distribution, this work utilized sampling techniques. A set of MDPs is sampled from the posterior distribution. An action is selected by considering the actions of all the sampled models and considering the expected gain and value of information in the uncertainty over the parameters.

In the original work, these models are resampled as frequently as possible. Followup work (Strens, 2000) suggested decreasing the resampling rate. By maintaining the set models for a longer period of time, the computation complexity is reduced. This work demonstrated that this is a viable strategy as one set of models can be used successfully to select a series of actions.

Bayes Adaptive MDPs (BA-MDPs) (Duff, 2002) presented a method to determine exact solutions for BRL. This approach formulated the problem as a partially observable Markov decision process. In the new POMDP, the states are a cross product of the original MDP states and the unknown parameters. Essentially, each state represented one possible parameter assignment. The observations were defined as the set of states from the original MDP. Given this formulation, POMDP planning techniques can be used to find a policy over the POMDP which allows the agent to consider the potential states and parameters which may be encountered. This provided a natural method for planning which balances exploitation and exploration and takes advantage of the BRL formulation. It has also been shown that the value function of the BA-MDP formulation takes a specific form, namely of a parameterized multivariate polynomial (Poupart et al., 2006). The BEETLE algorithm presented in this work is a modification of traditional POMDP solvers to take advantage of this property to speed up planning over BA-MDPs.

The size of the resulting model when determining exact solutions under BRL can be prohibitive. Recent work (Wang et al., 2005) addressed this issue with sparse sampling (Kocsis et al., 2001) to solve these larger models. Sparse sampling used the MDP as a generative model to generate samples to construct a lookahead tree,

which was then used to calculate a solution to the MDP. When issued a state-action pair, (s, a) , the generative model returned a sample resulting state, s' , and reward, r . This approach was able to find a solution to the model independent of the number of states, as well as capable of handling infinite models, thus providing a natural fit for the continuous models used in exact BRL methods.

The use of Bayesian reinforcement learning has also been demonstrated in a multiagent domain (Chalkiadakis & Boutilier, 2003). One of the primary questions for multiagent planning is considering the impact of one agent's actions on other agents and gauging their response. By taking this approach, the agent can account for its own actions, handle exploitation vs. exploration, and consider the influence on the decision-making of the other agents.

4.3 Bayesian Reinforcement Learning in Partially Observable Domains

The previous work described in section 4.2 assumes that learning is being conducted in fully observable domains. This implies that transitions between states are easily observed, and updating the necessary counts is straightforward. This problem becomes more difficult when working in partially observable domains. New techniques must be introduced to learn parameters for POMDPs as the underlying states are often not observed when a transition occurs.

The MEDUSA algorithm (Jaulmes et al., 2005) presented an approach for learning POMDP parameters. As with BRL for MDPs, a Dirichlet distribution was defined for each unknown transition probability. In the same fashion, distributions are also defined for each unknown observation probability. To address the issue of updates to the distribution, MEDUSA took advantage of the concept of an oracle,

based on ideas from active learning (Anderson & Moore, 2005). During execution, the learning could select to query an oracle for the true underlying state of the world. In this case, the oracle revealed (s, s') based on the previous action. Once this information was known, the Dirichlet parameters were updated as with the MDP formulation.

To determine the optimal policy based on the posterior, MEDUSA sampled a set of models from the posterior and calculated the policy for each sampled model. Actions were selected by examining the actions of each individual model and using a weight voting scheme. To keep this approximation of the optimal policy current, the models were periodically resampled. These models were also used to determine when to query the oracle for information. Heuristics based on the entropy of the actions of the models could be used to detect when significant confusion occurs, highlighting the need for information from the oracle.

In (Ross et al., 2007), BA-MDPs were extended to the partially observable domain. This work defined a Bayes Adaptive POMDP (BA-POMDP), where each state incorporates the states from the original POMDP as well as the Dirichlet distribution counts. As with BA-MDPs, this approach attempted to determine an exact solution over the models. While the state space in this model was infinite, bounds on the space of counts could be used to generate a finite, although still large, approximate model. The policy was determined through online planning methods (Paquet et al., 2005) for POMDPs. During execution, the agent performed a finite-depth lookahead in the BA-POMDP from the current belief. To maintain tractability, the belief state was approximated use Monte Carlo sampling methods.

Partially Observable BEETLE (Poupart & Vlassis, 2008) adapted the BEETLE algorithm for partially observable domains. PO-BEETLE also constructed a new model in a similar fashion to BA-POMDPs. This work demonstrated that the resulting belief state over this new model is defined by a mixture of Dirichlets, and the resulting value function represented as a linear combination of Dirichlets. By using this approach, PO-BEETLE was able to perform exact updates to the belief and value function.

A similar approach to BRL for POMDPs was seen in (Doshi et al., 2008). As with (Ross et al., 2007)’s BA-POMDP, a new POMDP was constructed combining the states of the original POMDP and the potential models. This work, however, sampled a set of candidate models from the distribution as with MEDUSA. Though in this case, an oracle returning policy information was used, as opposed to state information; policy information is often a more natural response from a operator working alongside an agent. One of the main contributions in this work was the use of Bayes risk for action selection. By considering the uncertainty over the sampled models, an action could be selected which minimizes the risk taken. This could also be used to help determine when to issue a query to the oracle for information. Updates to the distribution were performed using an EM approach, where the distribution over the states is first estimated, then used to update counts of the Dirichlet distribution.

4.4 Summary

Bayesian reinforcement learning has proven to be a rich mechanism for learning Markov models. This method for learning provides several advantages: a natural

mechanism for exploration versus exploitation, an efficient method for assimilating new data, and the ability to handle prior knowledge. Research has been conducted to develop BRL in fully observable as well as partially observable domains. These include analytic solutions, as well as methods to approximate the models and optimal policies.

However, there is still significant research to be done concerning Bayesian reinforcement learning. The research concerning partially observable domains is still limited. Exact methods are limited in the size of model which can be learned. Many algorithms assume the reward is an observable parameter which can prove to be a limitation in real-world environments. Finally, very few of the algorithms presented are demonstrated on complex domains. Experiments are limited to toy domains, or very small problems. In the following chapters, we will demonstrate a novel framework for learning POMDP parameters which takes advantage of Bayesian reinforcement learning and extends many of the ideas discussed here. Furthermore, we will demonstrate this framework on a real world interaction domain.

CHAPTER 5

A Framework for Online Parameter Learning

In the previous chapters, we have presented Markov models as a mechanism for representing stochastic environments. While decision-making using these models has been extensively studied, learning the model parameters still requires investigation. This is particularly true in partially observable environments. Applying traditional learning algorithms is difficult due to noise and incomplete information about the environment. We also presented Bayesian reinforcement learning as a technique for learning model parameters. While BRL has been extensively explored, there remain many open questions. Most research focuses on fully observable domains, while research in partially observable domains is still limited. Learning in partially observable environments is difficult due to the lack of complete information. In this chapter, we present a novel framework for learning in partially observable domains which builds on earlier work on Bayesian reinforcement learning. This approach centers around the use of policy information from an oracle to direct a gradient ascent update to the distributions.

5.1 Framework

As with most BRL approaches, the overall idea behind this framework is to represent the unknown parameters in the model as distributions instead of point estimates. The agent can then reason about the uncertainty of the parameters when performing action selection. This allows the agent to take actions which influence the

Algorithm 4 A framework for Bayesian reinforcement learning. In this algorithm, *ask_Oracle* and *resample_POMDP* are heuristics based on the domain. *Update_Distribution* is specific to the parameters being updated and is discussed in following chapters.

Initialize Beta/Dirichlet Distribution based on prior knowledge

Sample n POMDPs M_1, M_2, \dots, M_n from distributions

Solve M_1, M_2, \dots, M_n using any POMDP solution method

loop

if ask_Oracle() **then**

 Get action a_o from oracle

else

 Get action from models by voting

end if

 Execute Action

 Obtain Observation

 Update Belief States

for all M_i where $\pi_{M_i} == a_o$ **do**

 Update_Distribution()

end for

if resample_POMDP() **then**

 Remove k worst performing POMDP

 Sample k new POMDP from distribution

 Solve new models using POMDP solution method

end if

end loop

learning by considering the potential resulting information, as well as allowing the agent to select the optimal action given the variance of the parameters. Information gathered after action execution is used to update the parameter distributions through a learning update. The use of Bayesian reinforcement learning also allows domain knowledge to be used in the form of a prior distribution.

This framework extends BRL by presenting a novel approach based on the sampling used to approximate the policy. While previous work utilizes the sampling only for action selection, we take advantage of information from the samples to direct the learning. First, we utilize the information provided by the oracle. As part of the approximation of the policy, the framework has a set of sampled models. By comparing the policies of those models with the correct action provided by the oracle, the framework can select the models which are most representative of the environment. Second, once those models are selected, the distribution hyper-parameters can be updated based on the parameters of the correct models by utilizing a gradient ascent technique. These two steps are detailed in the following sections.

This approach provides several advantages. First, the framework is able to learn in partially observable environments. In fully observable environments, the updates to the distributions are straightforward as the state is observed. Instead, this framework relies on information from the oracle to help perform the updates in partially observable domains. Second, this framework does not require extending the underlying POMDP model. This approach avoids the problem of solving a large continuous POMDP as seen in other work. Additionally, this framework uses sampling to approximate the posterior distribution and perform action selection, reducing the

computation required for planning. Finally, by taking advantage of the oracle information, the framework can learn the reward parameter in domains where the reward is not actively observed. This aids in overcoming a drawback of many earlier Bayesian reinforcement learning systems which rely on the reward as an observed parameter to estimate the values or assume the reward is known.

The overall structure of our framework is described in Algorithm 3. There are four key steps. First, the system initializes a prior containing a rough model of the environment. Distributions defined by hyper-parameters which act as counts, such as the Dirichlet distribution or Beta distribution, make ideal choices of prior. Next, a set of n candidate models are sampled and an optimal policy is obtained corresponding to each. During the learning phase, an oracle is queried to know the most desired action. This query occurs if *ask_Oracle* is true. This action is then executed, and a posterior over the parameters is computed via *Update_Distribution*. The update to the distribution varies based on the parameters to be learned. This is addressed in the following chapters. Periodically, models are removed and new models are re-sampled. This rate of resampling is controlled by *resample_POMDP*. In this algorithm, *ask_Oracle* and *resample_POMDP* are domain-specific heuristics, and *Update_Distribution* is specific to the distribution used.

There are four key points to this framework: The use of a policy oracle, the representation of the parameters, the update to the distribution after data is acquired, and action selection using the representation. These are discussed in detail in the remainder of this chapter.

5.2 Policy Oracle

We propose an approach which relies on having an oracle which indicates the optimal action upon request. This optimal action is a reflection of the overall policy. Taking advantage of policy information from the oracle helps provide information about the environment to direct the learning. By having correct policy information, the framework is able to use this information while updating the distribution parameters. This is particularly important while working in partially observable domains. This oracle is helpful when learning observation and transition parameters, but critical in learning reward parameters which would otherwise be unobtainable as they often lack physical data.

Previous work on learning using an oracle (Jaulmes et al., 2005) relied on the oracle for full *state* information. While this is sometimes feasible, in many scenarios it is not realistic to have an oracle to provide such detailed information. Especially in the case of human-robot interaction tasks, we believe it is much more reasonable to request that the user indicate a suitable action at any point in time. Consider the case where the task requires joint human-robot control of the navigation functions. In cases where the robot is unsure of the user’s preferences, it is perfectly natural to ask the user to indicate the proper navigation action, yet it is much more difficult for the user to provide precise localization information. In contrast, in a task centered around verbal interaction between the robot and human user, it be quite natural to conceive of the oracle as a back-up operator, as many automated phone systems currently operate, which is available during an initial training period. With our particular wheelchair platform, there is always a human operator, in addition to the

primary user, present during training phases, to monitor performance of the robot. It is much easier for this operator to substitute for the robot by providing optimal actions whenever needed, than it would be to reveal the state, which can be a very abstract notion in the case of dialogue interactions.

5.3 Representation

For each unknown parameter, we assign a distribution defined by a set of hyper-parameters. In this work, we focus on the use of the conjugate family of distributions. In Bayesian reasoning, distributions are said to be *conjugate* if the prior distribution, $p(\theta)$, and posterior probability distribution, $p(\theta|x)$, are in the same family of distributions. In this case, the distribution of the likelihood, $p(x|\theta)$, is known as the *conjugate prior*. By considering conjugate families of distributions, we are able to reduce the complexity of representation by providing a closed-form update to the distribution parameters, as the updated posterior distribution from one iteration can act as the prior distribution in the next iteration. Of particular note are the Dirichlet and Beta distributions. The Dirichlet distribution is the conjugate prior of the multinomial distribution and thus is a natural fit for representing model parameters such as the observations and transitions. The Beta distribution is the conjugate prior of the binomial distribution, and can be viewed as the special cases of the Dirichlet and multinomial distributions when only two events are involved. Details of these distributions can be found in Appendix A.

These distributions act as priors over the model parameters and can be used to capture domain knowledge of a problem. If no information is known about a domain, we assume a uniform distribution over the prior. However, if there is some

existing domain knowledge, such as from expert knowledge, rough handcrafted models, or models learned from a similar, but different, problem, then the initial hyper-parameters can be set accordingly to reflect this information.

In this work, we assume rewards are within the range $[0, 1]$. We associate each $R(s, a)$ to be learned with a Beta distribution with two hyper-parameters, $\alpha_{s,a}$ and $\beta_{s,a}$, which act as a prior over a Bernoulli distribution.

When creating a new POMDP, M , we first sample a distribution for each $R_M(s, a)$ from the Beta distribution. We can then set $R_M(s, a)$ to be the expected value, p_M^{sa} , of the resulting Bernoulli distribution:

$$R_M(s, a) = p_M^{sa} \tag{5.1}$$

In this work, we assume observations are only dependant on the state, and thus focus on learning $O(s, z)$. However, it is straightforward to extend the same approach to observations dependent on the action as well, in the case of $O(s, a, z)$. To represent the observations for learning, we associate with each state, s , a Dirichlet distribution over the possible observations, z . Each of these distributions is defined by a set of hyper-parameters $\alpha(s, z)$, defining counts over the number of times z was observed in state s . This acts as a prior over a multinomial which becomes the observation distribution for the state. When a new POMDP, M , is created, a multinomial defining the distribution over the observations is sampled from the Dirichlet distribution for each state and action, providing $O_M(s, z)$ where

$$(O_M(s, z_1), O_M(s, z_2), \dots, O_M(s, z_q)) \sim Dir(\alpha(s, z_1), \alpha(s, z_2), \dots, \alpha(s, z_q)) \tag{5.2}$$

A similar representation is used for transitions for learning. We associate with each state, s , and action, a , a Dirichlet distribution over the resulting states, s' . Each of these distributions is defined by a set of hyper-parameters $\alpha(s, a, s')$, defining counts over the number of times the agent transitioned to s' after executing a in s . As with observations, these act as a prior over a multinomial which becomes the transition distribution for the state. When a new POMDP, M , is created, a multinomial defining the distribution over the transitions is sampled from the Dirichlet distribution for each state and action, providing $T_M(s, a, s')$ where

$$(T_M(s, a, s_1), T_M(s, a, s_2), \dots, T_M(s, a, s_n)) \sim Dir(\alpha(s, a, s_1), \alpha(s, a, s_2), \dots, \alpha(s, a, s_n)) \quad (5.3)$$

5.4 Distribution Update

In a fully observable domain, when an event occurs, the hyper-parameter reflecting that event is incremented. This is feasible because the underlying state is known. For example, when an agent transitions from state s to s' via action a , one might simply increment a hyper-parameter $\alpha(s, a, s')$. However, in a partially observable domain, the underlying state is not observed, thus making it more difficult to assign credit and perform a single update. Instead, we take advantage of the information provided by the oracle as well as the current models.

An update to the hyper-parameters occurs after the oracle has provided an optimal action, a_o . The hyper-parameters are updated in the direction of the distribution of each model which also selected the same action at the step. The intuition behind

this update is that the models which selected the correct action are a better representation of the environment, and thus the search should be moved in the direction of those models. Because this update occurs frequently, the influence of a model which is poor, but happened to select a correct action occasionally, will be less significant than a more accurate model which is correct more often.

To update the hyper-parameters of the Beta distribution over the rewards, we take advantage of the Bernoulli distribution sampled from the model. After querying the oracle, we perform the following update for each model, M_i , which agreed with the oracle:

$$\alpha(s, a) \leftarrow \alpha(s, a) + \lambda[(1 - p_M^{sa})] \quad (5.4)$$

$$\beta(s, a) \leftarrow \beta(s, a) + \lambda[p_M^{sa}] \quad (5.5)$$

where λ is the learning rate. This appears as the *Update_Distribution()* step in algorithm 3.

As with learning reward parameters, a direct update of the hyper-parameters is difficult due to the partially observable nature of the environment. A similar update is applied to adjust the Dirichlet hyper-parameters. After querying the oracle, we perform the following update for each model, M_i , which agreed with the oracle:

$$\alpha(s, z) \leftarrow \alpha(s, z) + \lambda Pr_M(z|s) \quad \forall s \in S, z \in Z \quad (5.6)$$

where λ is the learning rate. This appears as the *Update_Distribution()* step in algorithm 3.

A similar update can be applied to the hyper-parameters of the transitions. After querying the oracle, we perform the following update for each model, M_i , which agreed with the oracle:

$$\alpha(s, a, s') \leftarrow \alpha(s, a, s') + \lambda Pr_M(s'|s, a) \quad \forall s, s' \in S, a \in A \quad (5.7)$$

5.5 Action Selection

In this framework, action selection is not only necessary to move the agent towards the goal, but also serves as a mechanism to aid with the search through the parameter space. Because the parameters to be learned are represented as distributions, determining an action analytically can be difficult and expensive. As seen in chapter 4, one common approach to determining an exact policy is to create a new model where the states are the possible combinations of the states in the original model and the hyper-parameters values. While the policy can now be determined exactly, the resulting model is continuous and thus determining a solution difficult. Instead, the policy can be found by approximating the posterior. In this framework, models are sampled from the posterior distributions as representative examples to act as an approximation. Once the models are sampled, a policy is determined using one of the solution techniques described in chapter 3. Once a set of models is sampled and solved, a single action can be determined through a heuristic such as voting, where the action selected by the most models is chosen as the action to execute. This sampling approach helps control the computation required to determine the policy.

By increasing the number of samples, the necessary computation is increased as well as the quality of the approximation.

As the parameters change over time, these models must be periodically resampled to reflect current parameter distributions. The number of models resampled at each resample iteration represents a tradeoff. When a larger number of models are resampled, the computation time is increased as the policies must be regenerated each time. This may throw out models which have selected correct actions, potentially replacing them with poor models, and hurt the immediate performance of the system. Additionally, maintaining these models is advantageous because the update to the distribution is based on these well-performing models. However, resampling too few models may result in a very slow search through the POMDP space and thus require many iterations to locate appropriate models.

5.6 Example

As a illustrative example, assume we have a state, s , and action, a , with two observations, z_1 and z_2 . We represent the two unknown parameters, $O(s, z_1)$ and $O(s, z_2)$, as a Dirichlet distribution defined by two hyper-parameters, $\alpha_1 = 8$ and $\alpha_2 = 3$. We sample three models from the Dirichlet distribution with the values $O_1(s, z_1)$, $O_1(s, z_2)$, $O_2(s, z_1)$, $O_2(s, z_2)$, $O_3(s, z_1)$, $O_3(s, z_2)$. After finding the policy for each model, it is determined that models 1 and 2 have selected the same action at this step as the oracle. Thus, an update to the distribution hyper-parameters may occur as follows:

$$\alpha \leftarrow \alpha_1 + \lambda(O_1(s, z_1) + O_2(s, z_1)) \tag{5.8}$$

$$\beta \leftarrow \alpha_2 + \lambda(O_1(s, z_2) + O_2(s, z_2)) \quad (5.9)$$

where λ is a learning rate set to control the speed of the learning.

5.7 Discussion

This chapter presents a Bayesian reinforcement learning framework for learning parameters of partially observable Markov decision processes. The main idea behind this framework is to utilize information provided by the oracle to perform updates to the distribution over the parameters. This is necessary due to the partially observable environment. This framework uses sampling to approximate the posterior. These samples are used for action selection as well as to direct the gradient ascent in combination with the oracle information. This framework does not extend the underlying POMDP model, as is done in most related work. This helps to reduce the computational requirements.

This chapter addressed four main elements of the framework: a policy oracle which provides optimal action information to the agent when necessary, representation of the parameters using distributions, update to the distributions using information from the policy oracle, and action selection using the distributions.

This framework compares against the methods for Bayesian reinforcement learning in partially observable environments presented in section 4.3. BA-POMDPs (Ross et al., 2007) and PO-BEETLE (Poupart & Vlassis, 2008) are both exact methods for BRL. Both methods require a continuous POMDP to be constructed and solved. This could prove prohibitive for more complex domains. Furthermore, both of these

works center on learning the observation and transition parameters, but do not address the learning of the reward parameter. At this time, PO-BEETLE does not present empirical evidence.

The MEDUSA algorithm (Jaulmes et al., 2005) and the work by (Doshi et al., 2008) are the closest to this framework. These approaches both take advantage of an oracle and uses sampling to represent the posterior for action selection. MEDUSA assumes the oracle is able to provide state information to the agent. This simplifies the updates to the distribution, but such an oracle may be inappropriate or unrealistic in many domains. Additionally, MEDUSA only learns transition and observation parameters. The work by (Doshi et al., 2008) uses a policy oracle similar to this framework, but repeatedly suggests an action, which is accepted or rejected by the user. This work uses an expectation-maximization approach to update the distribution. However, this work was only validated on example POMDPs from literature.

The following chapters demonstrate instantiations of the framework for learning specific POMDP parameters, including the reward and observation parameters. These chapters will present the representations and updates used to learn these parameters, as well as results using experimental data.

CHAPTER 6

Learning Reward

In this chapter, we focus on the application of the framework presented in chapter 5 towards learning the reward parameters of partially observable Markov decision processes. Rewards help shape the policy by dictating the cost or payoff of performing an action. These are often the most difficult parameters to determine when crafting a model by hand as they are usually not grounded in physical data. Observation and transition parameters in the model can more often be determined from physical analysis, sensor models, and noise models. Furthermore, due to this often less tangible nature of rewards, many existing learning techniques cannot be applied to learning from data. While there are examples where rewards have a direct analogy from the physical environment, such as the cost in time or resources of a robot operating a sensor or moving a distance forward, this is not always the case. Intuitively, rewards attempt to capture a relationship between the actions, for example, allowing a robot to consider the cost of operating a sensor versus the risk of moving blindly through the environment.

In this chapter, we demonstrate the application of the framework to learning rewards. We will discuss the representation required to model the reward as well as the update based on the oracle information. Many learning methods, such as those discussed in section 3.3.2, assume the reward is an actual value returned by the environment which can be observed by the robot. This is sometimes not a realistic

assumption. Instead, we utilize the information from the oracle combined with the information from the sampled models to direct the learning. We demonstrate the ability of the framework to learn the reward parameters of the wheelchair intelligent interface presented in chapter 2. This represents a rich complex domain which can benefit from such a framework. We also present experiments to verify the impact of the main points of the framework, such as the use of the oracle and the use of sampled models to direct the learning. We will provide an in-depth empirical analysis of the properties of the framework.

6.1 Representing and Learning Reward

We begin by restating the representation and update used for the reward as presented in chapter 5. We focus on the case where the transition and sensor models are known, but information about the reward model is missing. Thus, $T(s, a, s')$ and $O(s, a, z)$, but not $R(s, a)$, are specified. The objective is for the agent to learn a reward function which results in a learned policy matching that of the oracle.

We first note that a linear transformation can be applied to the reward function of any MDP or POMDP without affecting the resulting policy. Thus to simplify the learning, assume all rewards are within the range $[0, 1]$. We associate each $R(s, a)$ with a Beta distribution with two hyper-parameters, $\alpha_{s,a}$ and $\beta_{s,a}$, which act as a prior over a Bernoulli distribution. When creating a new POMDP, M , we first sample a distribution for each $R_M(s, a)$ from the Beta distribution. We can then set $R_M(s, a)$ to be the expected value, p_M^{sa} , of the resulting Bernoulli distribution:

$$R_M(s, a) = p_M^{sa} \tag{6.1}$$

Because the environment is partially observable, and rewards are never truly observed by the agent, a probabilistic update of the Beta parameters must be performed. To do this, we take advantage of the Bernoulli distribution sampled from the model. After querying the oracle, for each model, M_i , which agreed with the oracle, we perform the following update:

$$\alpha(s, a) \leftarrow \alpha(s, a) + \lambda[(1 - p_M^{sa})] \quad (6.2)$$

$$\beta(s, a) \leftarrow \beta(s, a) + \lambda[p_M^{sa}] \quad (6.3)$$

where λ is the learning rate. This appears as the *Update_Distribution()* step in algorithm 3. The intuition behind this update is that models which agree with the oracle have captured some information about the world correctly and should have an influence on the direction of the learning. By moving the distribution towards these models, we can control the direction of search through the model space.

6.2 Experimental Methodology

The interaction manager presented in section 2 provides a domain which can take advantage of learning. The POMDP used in the experiments was largely handcrafted over time during the development phase of the project. This required multiple revisions of the model, repeated trials, and general adjustment of parameters until the desired behavior was obtained. A more efficient procedure would have been to provide a very rough initial model and then allow the system to tune itself automatically over time based on data from users. While observation and transition parameters

can be easily calculated from labelled data, determining the reward is not as straightforward. Adjusting the reward, for example, controls how often the system asks for clarification, determining how cautiously the robot proceeds. The purpose of these experiments is to apply the learning framework presented in the current chapter to learn the reward parameters for the interaction manager POMDP.

As noted earlier, the interaction manager POMDP contains 28 states, one for each potential intention of the user. There are 31 total actions, one action for each user intention, one general query action, and two specific query actions. The user intentions fall into one of two categories, movement and configuration, with each category having one associated specific query action. These experiments will be restricted to learning the costs associated with the specific query actions, $R(intention_a, query_a)$ and $R(intention_a, query_b)$, ie. the reward for asking an appropriate question, and the reward for asking an inappropriate question. By learning these rewards, the agent can determine how often to ask for more information, and when to use the specific requests versus the general request.

Parameter tying is applied for learning in this model. Parameter tying is a common technique in model learning which allows data to be used across several instances. When two parameters are tied, the value of those parameters is shared across both. This provides a speedup in learning as each datapoint example influences multiple parameters in the model. In this case, we assume the reward for an appropriate request is the same for both classes, and the reward for an inappropriate request is the same for both classes.

$$R(intention_a, query_a) == R(intention_b, query_b)$$

$$R(\textit{intention}_a, \textit{query}_b) == R(\textit{intention}_b, \textit{query}_a)$$

To facilitate experimentation, an environment was set up to replicate the user experiments conducted in the rehabilitation center. It is assumed the handcrafted POMDP is correct and, thus, the actions selected using its policy are optimal for the domain. This POMDP essentially takes the role of what would normally be the user in a deployed system. Two different procedures are used to generate the observations throughout the experiments. These are used to demonstrate the framework in different types of domains.

The first procedure, referred to as *data-based*, uses data collected from real users as observations. Phrases transcribed by the speech recognizer during user tests are sorted and labelled based on users' true intention. This provides multiple examples to express each intention, such as "move forward," "drive forward," and "drive the wheelchair forward," as examples of "move forward." Because these phrases are post-speech recognition, they also include misrecognitions such as "move ford" for "move forward" and "turn less" for "turn left." Each intention in the resulting data set has between 5 and 30 examples. This procedure for generating observations has the benefits of validating the algorithm on real world data collected from users. However, it is sensitive to the quality of the speech recognizer and the size of the data set collected.

During experimentation, when given an action, a , a random user intention is selected from the list of possible intentions based on the known transition function.

Next, an example phrase of that intention is selected at random. This phrase is passed to the grammar module as done with the interface experiments, which relays a parsed output to the interaction manager. Finally, the interaction manager updates the belief state based on the parse and selects an appropriate action based on the optimal policy. These resulting belief state and action are used by the learning algorithm to update the models. Because we assume the observation model is known, this belief can be shared across all the models, with the action acting as the oracle action.

The second procedure, *generative-based*, uses the underlying interaction manager POMDP itself as a generative model. When given an action, a , the POMDP transitions from the current state, s , by selecting s' based on the transition model, $T(s, a, s')$. It then draws an observation, z , based on the observation model $Z(s', a, z)$, and updates the belief state accordingly. Essentially, the POMDP is being used as a primitive simulator of the environment. This procedure is less sensitive to the data set, but may not be fully representative of the deployed domain. It does, however, serve to demonstrate some of the properties of the learning framework. Removing the dependency on the dialogue data also demonstrates that this framework can be applied on any domain which can be modeled as a POMDP.

For these experiments, we assume the oracle is queried at every step. Approaches to selecting queries are addressed in the related work in section 4.3. For these experiments, we focus strictly on the impact of the learning mechanism. Reducing the number of queries to the oracle simply increases the time for learning.

Throughout these experiments, 20 POMDPs are sampled and maintained from the Beta distribution. Planning on the sampled POMDPs is executed for 5 minutes. Learning is conducted over 30,000 steps using the *data-based* or *generative-based* procedures described above. A resample iteration occurs every 1000 steps. Unless otherwise noted, the three worst-performing POMDPs are removed, and three new POMDPs are sampled from the current Beta distribution. Experiments were repeated five times, with the average results reported. Because the reward functions cannot be directly compared, the models must be judged based on the performance of the associated policies. During execution, the percentage of time each model selected the same action as the oracle since the last resample iteration is recorded. This is reported as the performance of the models. For these experiments, a uniform prior is used where the hyper-parameters of the Beta distribution are initialized to 1, thus assuming no a priori knowledge of the environment.

6.3 Results

Figure 6–1 shows the performance of the learned models over time using data from the *data-driven* procedure. This figure reports the average performance and standard deviation over the performance of each of the models, where performance is measured as the percentage of time the models were able to select the same action as the oracle since the previous resample iteration. Because only two parameters are being learned, the models are initially able to select the correct actions 87% of the time. Initially, there is significant variation in the performance of the models, as demonstrated by the bars in the figure. As more data is accumulated, we see the average performance of the models improving. Additionally, the standard deviation

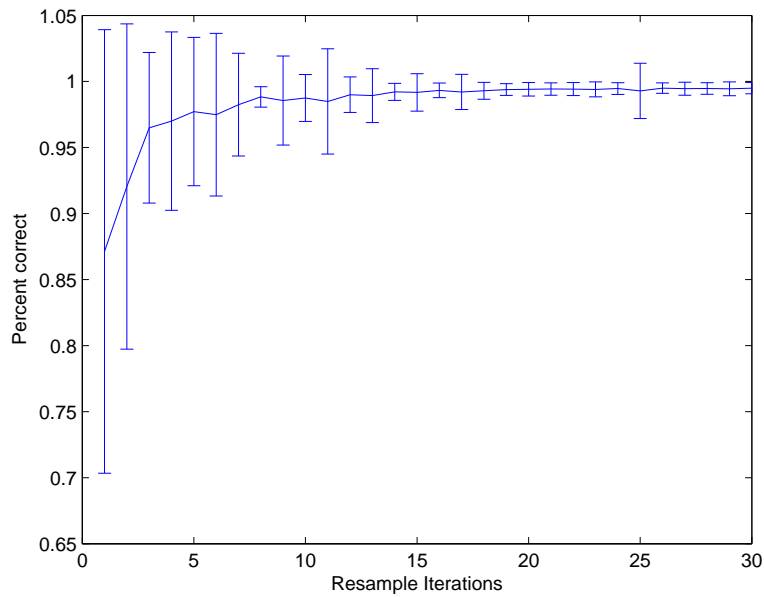


Figure 6–1: Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. The graph presents the average and standard deviation of the performance over all the models. Observations are generated using a *data-based* procedure.

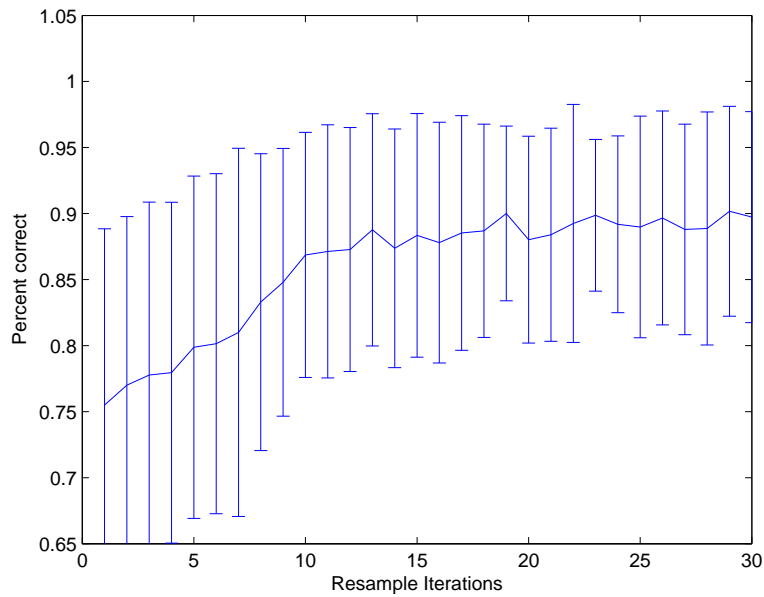


Figure 6–2: Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. The graph presents the average and standard deviation of the performance over all the models. Observations are generated using a *generative-based* procedure.

of the models is also reduced. This is a consequence of the improvement in the parameterization of the distribution. With more data, the hyper-parameters become a better representation of the environment, and generate more accurate models, and thus fewer poorly performing models. By the 15th resample iteration, the models are on average 99% correct with little variation in the performance of the individual models. Initially, there is a rapid increase in the performance as the sampling removes the most outlier models. As learning slows, the distribution update continues to improve the model qualities.

A similar experiment was executed using the *generative-based* procedure. These results are shown in figure 6–2. Again, we notice an improvement in the models as the number of iterations increases. Because the *data-based* procedure uses recognized phrases to generate the observations, many more belief states are visited, providing more data for the learning framework to generalize across. The *generative-based* procedure only works with the set of observations defined by the POMDP, limiting the number of explored belief states. This results in a slower overall learning and higher variance in the performance.

The number of models resampled at each resample iteration represents a tradeoff. As more models are sampled, the framework is able to move more quickly through the POMDP space. However, by keeping more well-performing models, the framework is able to select better actions at the current moment. Maintaining better models also helps update the distribution rapidly as the update is based on the models which agree with the oracle. Figure 6–3 demonstrates the impact of selecting the number of POMDP models resampled at each resample step. Note that the scale is reduced

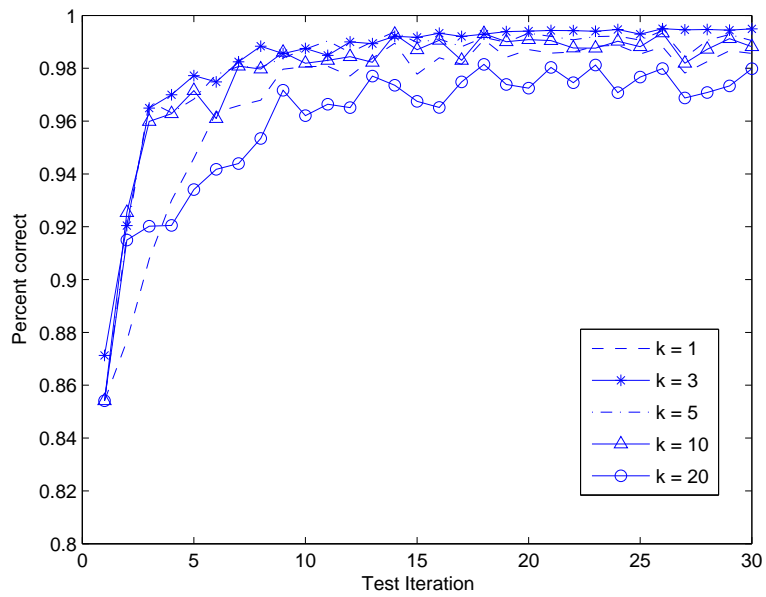


Figure 6–3: Average performance of sampled POMDP models over time with varying number of models resampled at each iteration. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. Observations are generated using a *data-based* procedure.

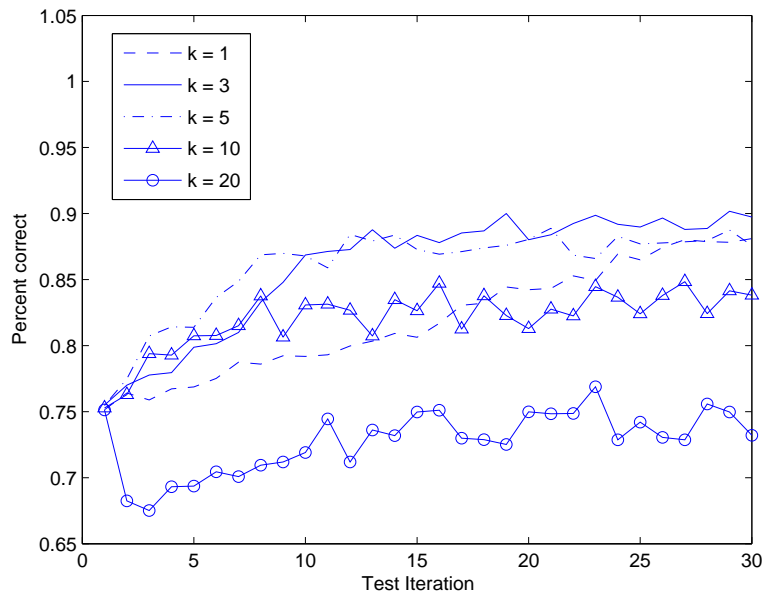


Figure 6–4: Average performance of sampled POMDP models over time with varying number of models resampled at each iteration. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. Observations are generated using a *generative-based* procedure.

for clarity. For small numbers, the learning occurs more slowly as the framework takes many iterations to remove the poor initial models. When only one model is resampled at each iteration, learning is very slow, although does converge to the proper value. For large values, such as 10 and 20, the system also performs less well as better performing models are removed. With such a high turnover in models, the performance of the system also becomes highly variant.

Of particular note is the instance where all 20 models are resampled at each iteration. The figure demonstrates that the performance of the system does improve over time. This reflects the impact of the update to the distribution hyper-parameters. Because all of the models are being replaced, the performance improvement in this case is due primarily to the resampling of models. While the resampling is a powerful factor for learning, particularly early in the process, the hyper-parameter update guides the direction of the search. If this was not the case, there would be no improvement with more data, as demonstrated in the following set of experiments. Performing the same experiment using the *generative-based* procedure exhibits a similar trend, as show in figure 6–4. As with the earlier experiment, the learning is, however, slower due to the reduction in the variation of the data set.

The results in figures 6–3 and 6–4 suggest that the model resampling is an important aspect of the learning process. To explore this further, experiments were conducted where no update to the distribution was made over time. Learning is allowed only via resampling to remove and replace the worst performing models. Models at each iteration are sampled from the initial prior. As expected, the system does improve some as the worst models are simply weeded away. However, the system

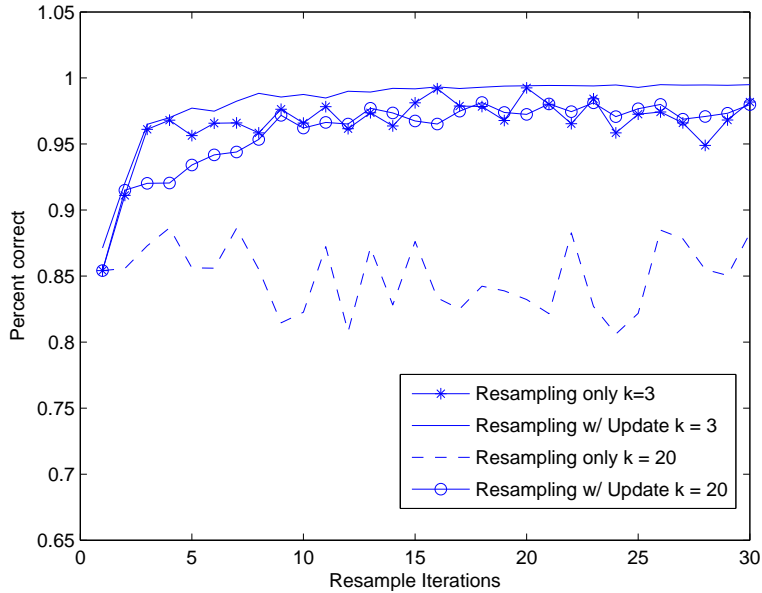


Figure 6-5: Results of learning with only resampling and no update to the distribution using a *data-based* procedure.

still demonstrates high variance in the model performance. These results are shown in figure 6-5. Results from figure 6-3 are plotted for comparison against resampling while performing the update.

We now examine the impact of limiting which models contribute to the distribution update. As stated earlier, the hyperparameters are updated only using information from models which selected the same action as the oracle action. Figure 6-6 shows experimental results where all the models, regardless of actions selected, were used to update the distribution at each step. Experiments were conducted with resampling of three models, and resampling of all 20 models. Results from experiments conducted using the standard update rule are presented for comparison. For

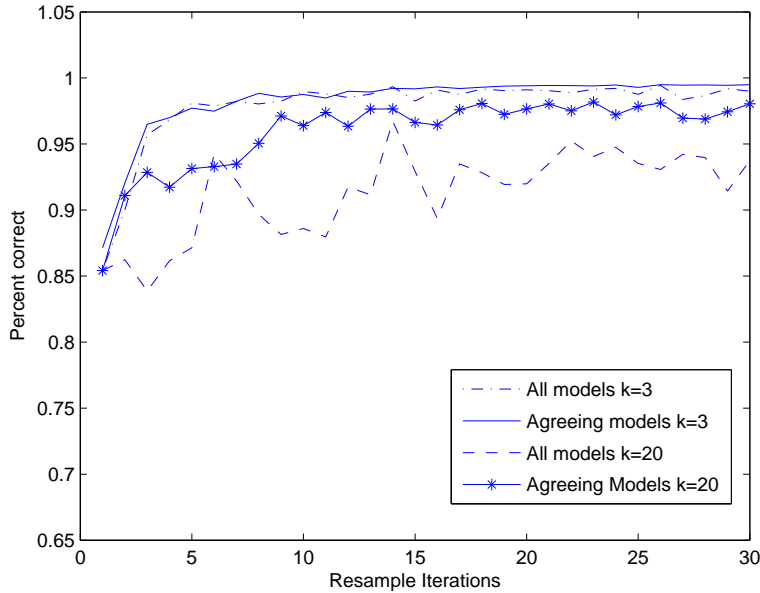


Figure 6–6: Results of learning on *data-based* approach when the distribution update is not limited to the models which selected the correct action.

the case with three models resampled at each resample iteration, learning is still occurring, with the system converging to almost the same point as in the earlier experiments. The impact is clear once the influence of resampling away poor models is removed. In the case where all 20 models are replaced at each iteration, it is clear that updating using all models reduces the rate of learning significantly.

Finally, we examine the influence of priors on the learning. Figure 6–7 presents the results of experiment conducted by varying the initial priors. In this case, the counts of the hyper-parameters for the informed prior are generated based on groundtruth values. These initial counts are approximated based on the correct distribution. This represents the case where information about the domain is known in advance and can be used to provide a starting point. For the uniform prior, where

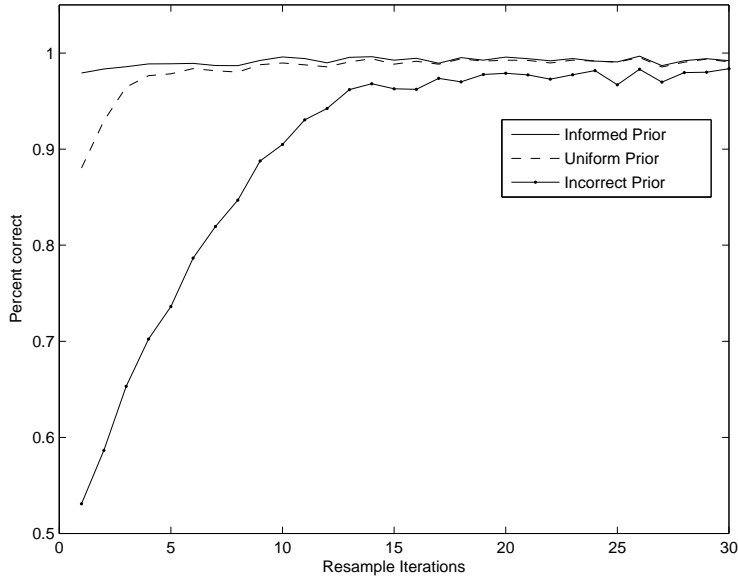


Figure 6–7: Learning of reward parameters given an informed, uniform, and incorrect prior. Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. Observations are generated using a *data-based* procedure.

no information is known, all initial counts for the Dirichlet distribution are set equal. The incorrect prior is determined by informed prior counts and reversing the counts, ie. $\alpha(s_0, a_0) = \beta(s_0, a_0)$, $\beta(s_0, a_0) = \alpha(s_0, a_1)$, ... In the case of the informed prior, the initial models perform well demonstrates slight improvement. With the incorrect prior, initial performance is poor. However, the framework is able to adjust the models over time. Learning is slower as the framework is sampling from a poor area of the parameter space.

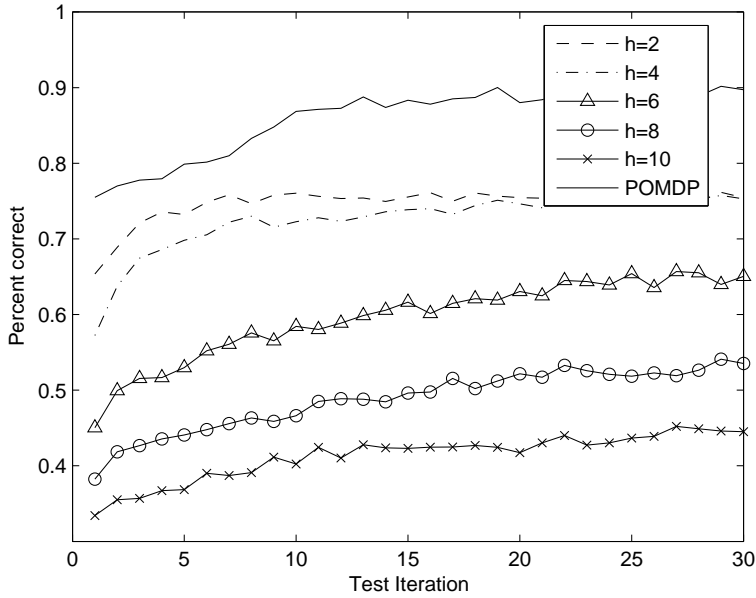


Figure 6-8: Results using a supervised learning approach.

As a baseline for comparison, an example-based supervised learning approach (an implementation of instance-based learning (Mitchell, 1997)) was also implemented. The goal of this experiment is to examine a model-free approach to determining the policy from experience. This method uses a *generative-based* approach to generate the data. During the learning phase, the history, H of actions and observations is recorded as a set of action-observations pairs, (a, z) :

$$H = a_0, z_1, a_1, z_2, \dots, a_{t-1}, z_t \quad (6.4)$$

During evaluation, the supervised learning method maintains a recent history, H_r , of length h :

$$H_r = a_{t-h}, z_{t-h}, \dots, a_{t-1}, z_t \quad (6.5)$$

At each step, an action is selected by locating the sequence in the learned history that most closely matches the recent history window and selecting the action taken at the next step. To match the sequence, H is searched for the sequence with the most elements that match H_r in the same order. If a sequence is found, the action immediately following the sequence is selected. If no matching sequence is found, a random action is selected. Similarly, if multiple examples of the sequence are found, with different actions selected, one instance is selected at random. Experiments were conducted for different recent history lengths, h . This supervised learning method represents a model-free method to learn the policy directly from the actions and observations recorded. This is similar to early instance-based methods for reinforcement learning on POMDPs such as NSM (McCallum, 1994) or USM (McCallum, 1995). However, those methods assume the reward is observed, and thus are able to maintain estimates of the values with the instances. We do not assume the reward is observable, and thus must directly work with the action selected.

To test the supervised learning approach, a simple evaluation phase is run periodically. Beginning with the initial belief state, at every step, an action is selected by the oracle, the action is executed, an observation is returned, and the belief state is updated. The action which would have been selected at each step by the SL approach is compared to the true action. The number on which they agree is tallied. For our experiments, this simulation is run every 1000 iterations to correspond with the POMDP resampling iterations, and is run for 1000 steps.

Figure 6–8 shows the results of learning on this interaction manager using our Bayesian method (POMDP), as well as the supervised learning method (SL) for different values of h . The results show that learning is occurring in both methods. For small values of h , the supervised learning improves very quickly as the history is built up. However, it never approaches the performance of the POMDPs. As h increases, the performance degrades quickly as sequences occur which are not encountered before, or sequences occur with multiple examples of actions chosen.

6.4 Related Work

Research has focused on various methods to provide guidance to agents for learning. Other work on *oracles* (Armstrong-Crews & Veloso, 2007) allows an agent to actively request information when deemed necessary. Oracles are often combined with the idea of *active learning* (MacKay, 1992; Anderson & Moore, 2005) which helps measure the value of information to be gained by requesting help. These approaches for learning are related to *imitation* or *apprenticeship learning* (Atkeson & Schaal, 1997; Price & Boutilier, 1999; Bagnell & Schneider, 2001; Lockerd & Breazeal, 2004) where an agent observes another agent or human perform a task, and must then replicate the task.

To handle more general MDPs, inverse reinforcement learning (Ng & Russell, 2000) (Abbeel & Ng, 2004) was introduced as a method for learning a reward function of an MDP which explains an observed policy. This technique formulates the MDP as a linear program and finds the reward value which maximizes the value function, using the observed policy as constraints. A very similar approach is seen in work on maximum margin planning (Ratliff et al., 2006). While these were able to determine a

reward function, heuristics had to be used to prevent degenerate cases. For example, if every action returns 0 reward, the policy under this reward structure agrees with the optimal policy. To handle these issues, Bayesian inverse reinforcement learning (Ramachandran & Amir, 2007) was introduced. In this work, a prior over the reward function is used to direct the solution. Inverse reinforcement learning has also been demonstrated on POMDPs (Choi & Kim, 2009). This is performed by formulating a linear program over the POMDP and determining the rewards which maximize the value function. While inverse reinforcement learning is able to determine a reward, a key step in these approaches is solving a linear program. This is computationally intensive and is only demonstrated on small toy domains.

6.5 Conclusion

The experiments presented in this chapter demonstrate the various aspects of learning using our Bayesian reinforcement learning framework. This framework achieves learning through two major mechanisms: resampling of reward models and updating the distribution over the reward function. Through the experiments, we have explored the learning capabilities as the parameters of these aspects are varied.

The resampling generally proves to be a powerful mechanism in learning. This idea of resampling is related to Monte Carlo search methods which utilize repeated random sampling and evaluation to search a space. These approaches are useful in domains where analytical approaches may be difficult or impractical. Because the metric for removal in our framework is based on model performance, simply removing poorly performing models and maintaining high performing models produces a system that performs well overall. The resampling can be seen as moving the search

through the space of potential models. The more models are selected for removal, the larger the step through model space. However, as the experiments also show, moving too quickly may miss the optimal set of models as well as reducing performance at the current moment.

The distribution update serves to control the direction of search through the space by dictating which new models are sampled. While the experiments show the resampling having a larger impact on the speed of search, they also show that updating the distribution parameters contributes to the search. Furthermore, the method for updating the distribution is also significant. Because rewards often do not have a real world basis, such as with observations, the rewards must be inferred through a different mechanism. In this case, we rely on models which select the same action as the oracle to suggest the direction of search based on the idea that they have captured some of the nature of the correct models.

When comparing a model-free supervised-learning approach against our model-based BRL approach, we find that BRL perform significantly better. Part of the reason for this is the information encoded into the transition model which is lost in the model-free methods. For the interaction manager used here, looking at a very short history is sufficient for determining the next action. However, as the model becomes more complex, more history information is needed to determine context. Because the supervised learning method must be explicitly told the length of history to consider, it currently has no way to determine the optimal length. Our POMDP approach, however, maintains history automatically through the use of the belief states and is not affected by this problem.

This is particularly crucial in the interaction domain because the dialogue system essentially resets itself after each completed action. It is assumed that after one task is complete, the user can request any task. This information is easily encoded into the POMDP through the transition function. A typical dialogue may begin with the user issuing a “turn left” command. The system issues a query, as “turn left” is easily confused with “turn right.” The user issues another “turn left” command at which time the system commits to turning left. In another case, the user issues a “stop” command which the system immediately identifies as the correct action, as no other action sounds similar. In the first case, the dialogue involved two interactions. In the second case, only one. The POMDPs approach allows the relation of the two phrases, due to the information in the transition models. However, the supervised learning approach assumes a fixed history length. Examining a history which is too short may not capture enough of the context to make the correct decision, while a history which is too long captures too much information, looking into a previous iteration which is not relevant to the current action. While this can be handled using domain-specific heuristics, the model-based system is able to capture this information implicitly.

There are other advantages to explicitly maintaining a model as opposed to using a supervised learning strategy. Models can be transferred between users, either by starting with the model itself or converting the model into priors and restarting the learning process. An interaction manager trained on one user can be transferred to another user. While this new model may not necessarily be tuned for the new user, it may be much closer than a rough hand crafted initial model, especially for elements of the problem which are not user dependant such as noise from the speech

recognition. Finally, the supervised learning method only provides information about experiences which have been explicitly encountered as direct matches need to be found in the history. By using models, generalization can occur more easily to novel experiences. POMDPs maintain a belief state, which encodes the information from the transition and observation models. These belief states are continuous, so belief states which are never explicitly encountered still have policy information as a result of the way policies are maintained. While more sophisticated versions of the supervised-learning approach could be developed, they would still lack the benefits of the model-based approaches.

The results in this chapter focus on learning the reward parameter of a POMDP, particularly centered around the wheelchair interaction manager domain. In the following chapter, we will address learning additional POMDP parameters, demonstrate the framework on additional domains, and compare the framework against additional approaches for learning.

CHAPTER 7

Learning Observations

In this chapter, we discuss the extension of the Bayesian reinforcement learning framework presented in chapter 5 to observation parameters, in addition to reward parameters. The need to use observations to determine the agent's state is one of the defining elements of partially observable environments. Observations represent an agent's perception of the environment and capture the error from this perception. Sensors on robots are often noisy, returning false positives or incorrect distances. In cases where observations are based on a secondary technique, such as calculations from a vision or speech recognition system, the error rate of that technique must also be considered. Even if sensors were able to perfectly perceive the immediate surroundings, they often could not provide a complete view of the world. Ambiguity arises from situations with perceptual aliasing, where different states of the environment appear identical such as with a robot in front of one of many doors in a hallway. A similar case occurs using speech recognition system where the phrase may not provide sufficient information to correctly determine the intention of the speaker even if the phrase was recognized correctly.

In light of these types of errors, proper models of the observation emission probabilities are necessary for the agent to perform action selection in these types of environments. While observations are sometimes grounded on physical characteristics of a system that can be modeled, an alternate procedure would allow the system

to learn the parameters automatically. This allows for a more general approach which can be applied to a variety of domains. As with learning of the reward parameters, this would also reduce the amount of domain knowledge required and help determine model parameters which may be difficult for the user to generate by hand.

In this chapter, we demonstrate the application of the Bayesian reinforcement learning framework presented in chapter 5 to learning observation. We present the representation of the parameters and the associated update. As with the rewards, we will validate the framework on the data collected for the wheelchair intelligent interface. We also present results combining learning of observations and rewards concurrently.

7.1 Representing and Learning Observations

We focus on the problem of learning the observation model assuming the reward and transition models are known. Thus, $T(s, a, s')$ and $R(s, a)$ are known, but not $O(s, z)$. The objective is to determine the observation parameters which replicate the policy of the oracle. It is important to note that these may not be the same observation parameters as used by the oracle.

We associate with each state, s , a Dirichlet distribution over the possible observations, z . Each of these distributions will be defined by a set of hyper-parameters $\alpha(s, z)$, defining counts over the number of times z was observed in state s . This acts as a prior over a multinomial which becomes the observation distribution for the state. When a new POMDP, M , is created, a multinomial defining the distribution over the observations is sampled from the Dirichlet distribution for each state and

action, providing $O_M(s, z)$ where

$$(O_M(s, z_1), O_M(s, z_2), \dots, O_M(s, z_q)) \sim \text{Dir}(\alpha(s, z_1), \alpha(s, z_2), \dots, \alpha(s, z_q)) \quad (7.1)$$

Essentially, the sampled observation distribution is based on the counts determined from data. Details of Dirichlet distributions can be found in Appendix A.

As with learning reward parameters, a direct update of the hyper-parameters is difficult due to the partially observable nature of the environment. A similar update is applied to adjust the Dirichlet hyper-parameters. After querying the oracle, we perform the following update for each model, M_i , which agreed with the oracle:

$$\alpha(s, z) \leftarrow \alpha(s, z) + \lambda \text{Pr}_M(z|s) \quad \forall s \in S, z \in Z \quad (7.2)$$

where λ is the learning rate. This appears as the *Update_Distribution()* step in algorithm 3.

7.2 Learning Rewards and Observations

We now combine the mechanisms for learning rewards and observations to determine both sets of parameters concurrently. In this case, we assume that $T(s, a, s')$ is known, but not $R(s, a)$ or $O(s, z)$. The objective is to learn the reward and observation functions which replicate the policy of the oracle. To achieve this, we associate a Beta distribution with each $R(s, a)$ with parameters $\alpha(s, a)$ and $\beta(s, a)$ as seen in section 6.1, and a Dirichlet distribution for each state, s , over the observations, z , with hyperparameters $\alpha(s, z_1), \alpha(s, z_2), \dots, \alpha(s, z_q)$ as seen in section 7.1.

7.3 Evaluation

As with the rewards, the observations used in the interaction manager in section 2 were determined by handcrafting and multiple revisions. While the resulting model performs well for the task, the construction required extensive domain knowledge and multiple iterations. By extending the Bayesian reinforcement learning framework, the observation parameters can be learned automatically from data. This technique is general to any problem which can be formulated as a POMDP, reduces the burden on the system designer, and allows the system to adapt over time. By learning the observations and rewards concurrently, we can further reduce the need for highly accurate handcrafted models. In this section, we validate the extensions to the framework by demonstrating learning of observation parameters on the interaction manager domain presented earlier. We also bring together the work from chapter 6 and demonstrate learning of reward and observation parameters concurrently.

Experiments were conducted to evaluate the learning of the observation parameters in a similar fashion as in section 6.2. For these experiments, we focus on the *generative-based* method for experiments, as the *data-based* method utilizes a separate domain-specific technique for updating the belief state using the input.

As before, the handcrafted POMDP acts as the oracle. The policy for the POMDP is determined using a point-based approximation method, and the belief state maintained during execution. When the oracle is queried for the policy information, the action for the current belief state of the ground truth POMDP is returned. In a deployed system, these actions would be provided by a user. Using the model as an oracle allows us to run repeated simulation experiments rapidly.

Throughout these experiments, 20 POMDPs are sampled and maintained from the Dirichlet distribution. Planning on the sampled POMDPs is executed for 5 minutes. Learning is conducted over 30,000 steps using the *generative-based* procedures. A resample iteration occurs every 1000 steps. The three worst-performing POMDPs are removed, and three new POMDPs are sampled from the current Dirichlet distribution. Experiments were repeated five times, with the average results reported. Because the observation functions cannot be directly compared, the models must be judged based on the performance of the associated policies. During execution, the percentage of time each model selected the same action as the oracle since the last resample iteration is recorded.

Parameter tying is again used, such that all observations which have nearly identical distributions are tied to the same counts. This reduces the number of parameters to learn from thousands to approximately 20. We focus on learning the 8 parameters which define the most common distributions, ignoring the observation distributions which are almost deterministic. Furthermore, in these experiments a prior is used. The prior in this experiment is very roughly approximated by selecting initial counts in proportion to the ground truth distribution. This approximates a situation where a rough model of the world is known.

To provide a comparison, the Baum-Welch algorithm (Rabiner, 1990) was implemented and tested on the model. The Baum-Welch algorithm is an expectation-maximization algorithm (Dempster et al., 1977) traditionally used for learning parameters in hidden Markov models (HMMs) given only a sequence of observations. Hidden Markov models are similar to POMDPs in that they have an unobserved

underlying state which must be inferred through a series of observations. The main difference being that HMMs are used for monitoring processes without decision making. They can be thought of as POMDPs with a single action. The algorithm works in two stages. First, the probability of each state is calculated based on fixed models parameters and the observation sequence. Second, the model parameters are adjusted based on these probabilities. This iteration is repeated several times to improve the models. Baum-Welch has become a popular method for training HMMs and is used extensively in many applications. However, it tends to require a large amount of data for training, especially as the models become more complex. For these experiments, the parameters of the model are trained using data sets of size 1000, 5000, 10000, and 30000. Five iterations of training are used, with a uniform prior over the parameters. The policy for the resulting POMDP is determined. For evaluation, 1000 steps of the ground truth POMDP are simulated. The resulting action of Baum-Welch-trained model policy is compared against the oracle policy at each step. The results report the number of times the Baum-Welch model select the same action as the oracle.

Figure 7–1 shows the average and standard deviation of the performance of the learned models over time using the Bayesian reinforcement learning framework. In this case, because only some of the parameters are being learned and a prior is used, the models perform reasonably well initially but still select an incorrect action often. There is also considerable variance in the actions selected by the models, representing significant disagreement among the models. As more data is collected, the performance of the models improves and the variance in the selected actions decreases.

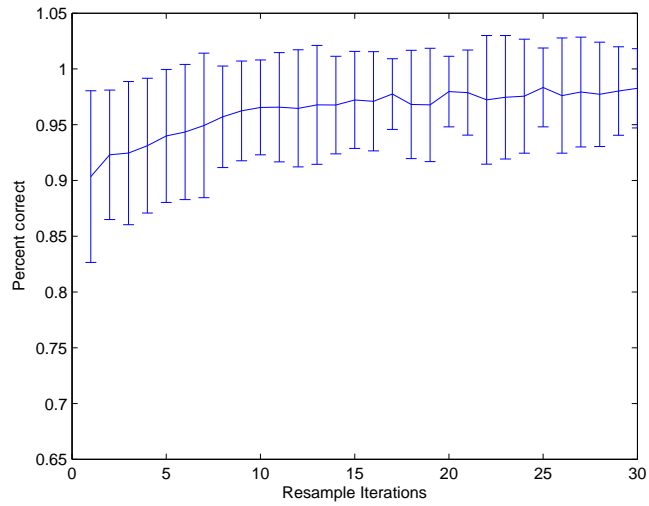


Figure 7–1: Learning Observations - Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. The graph presents the average and standard deviation of the performance over all the models. Observations are generated using a *generative-based* procedure.

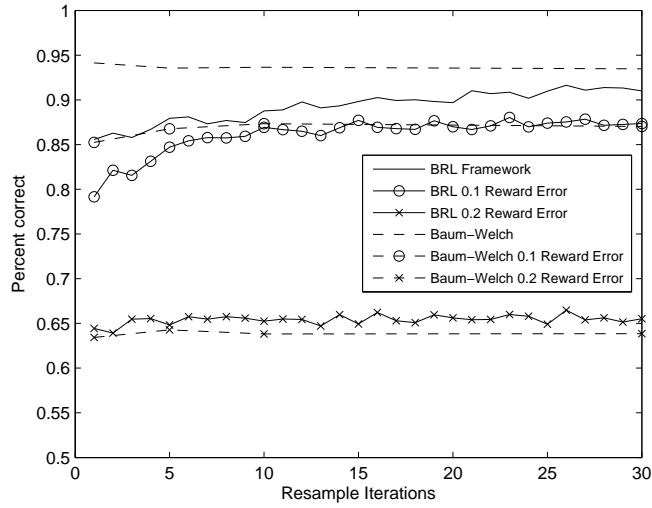


Figure 7-2: Average performance of models of the dialogue domain trained using Bayesian learning framework and models trained using Baum-Welch given error in the reward function. Using BRL, the models are resampled every 1000 steps. Using BW, models are retrained after every 1000 data samples collected. At each resample/retraining iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded.

This is similar to the behavior observed while learning the reward parameters. However, the learning occurs more slowly in this case as the observation space has higher dimensionality and more parameters are being learned.

The results for Baum-Welch are shown in figure 7-2. In the ideal case, the model learned by this method is able to chose the correct action 93% of the time. This may be due to information lost during the parameter tying. The results show that the Bayesian reinforcement learning framework approaches the performance of the Baum-Welch approach. However, POMDPs trained using Baum-Welch assume the reward function is correctly known. This approach does not contain any mechanisms for learning the reward parameters or compensating when the reward

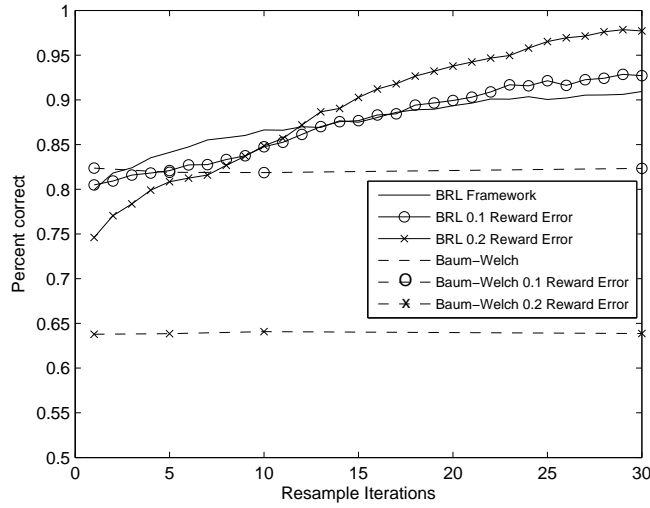


Figure 7–3: Average performance of models of the tiger domain trained using Bayesian learning framework and models trained using Baum-Welch given error in the reward function. Using BRL, the models are resampled every 1000 steps. Using BW, models are retrained after every 1000 data samples collected. At each resample/retraining iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded.

parameters are incorrect. To demonstrate this, experiments were conducted where error is introduced into the reward parameters. Figure 7–2 also shows results with errors of 0.1 and 0.2 (where the reward is in the range $[0, 1]$). When the reward is not correctly known, the Bayesian reinforcement learning framework is able to select more actions correctly than the models trained with Baum-Welch. This becomes more significant in the case of the larger error.

To further demonstrate the advantage of learning while considering the policy, the experiment was repeated on the classic tiger POMDP from literature, with 2 states, 2 observations, and 3 actions. The reward for the *listen* command was adjusted by 0.1 and 0.2, and a uniform prior is used for the Bayesian framework. These

results are presented in figure 7–3. In this example, the Baum-Welch algorithm is able to estimate the correct observation parameters associated with the original model and the resulting POMDP is able to select the correct actions in the case of the correct rewards. However, as observed in the earlier example, when error is introduced into the reward, the policy for the POMDP becomes increasingly incorrect. By using the Bayesian reinforcement learning framework, which considers the policy when adjusting the model parameters, models can be found which more accurately reflect the policy of the oracle. It is worth noting the behavior of the Bayesian framework when error of 0.2 is introduced. Under this reward, the framework is able to learn more quickly, and eventually learn a better model, than the cases with less or no error. In this case, the resulting policy for this reward is less sensitive to the observations, thus when sampling, a wider range of observations can result in matching policy. When the reward had less error, the observations required for the resulting policy to match the oracle had to fall within a smaller range.

Figure 7–4 presents the results of using the Bayesian reinforcement learning framework to learn both the reward parameters and the observation parameters concurrently. For these experiments, Beta distributions are used to represent the reward parameters, as seen in section 6.2 and Dirichlet distributions are used to represent the observations as seen in section 7.1. These results demonstrate a similar trend as observed learning each parameter individually. As more data is received, the models are better able to select the correct actions. The variance in the model performance also decreases.

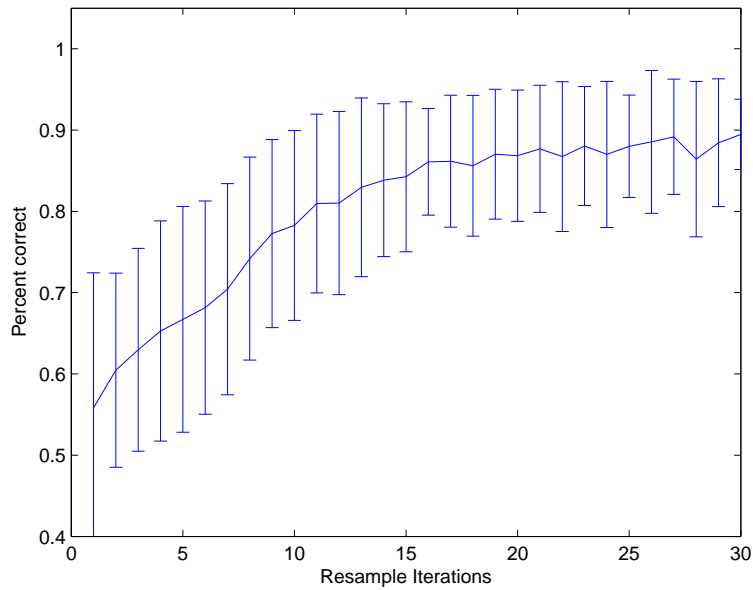


Figure 7–4: Learning Rewards and Observations in Dialogue domain - Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. The graph presents the average and standard deviation of the performance over all the models. Observations are generated using a *generative-based* procedure.

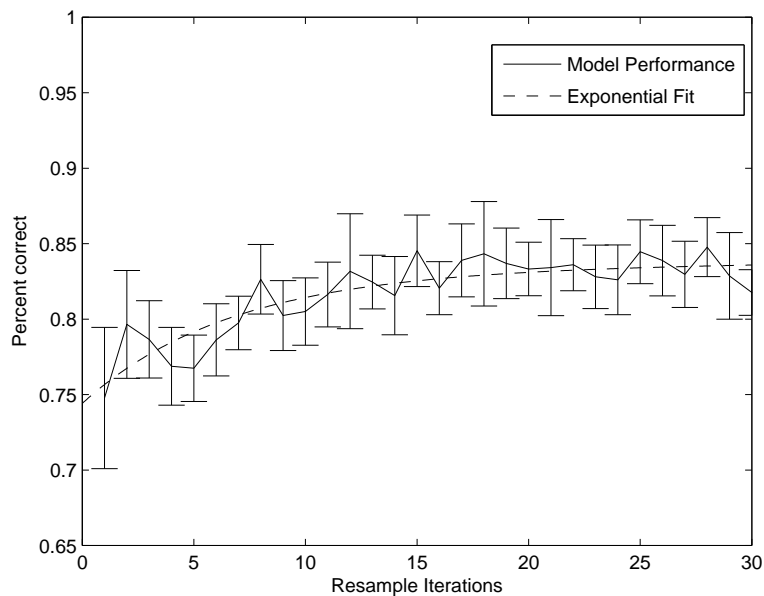


Figure 7–5: Learning Rewards and Observations on SACTI model - Average performance of sampled POMDP models over time. A resample iteration occurs every 1000 observations. At each resample iteration, the percentage of times each model selected the same action as the oracle over the past 1000 steps is recorded. The graph presents the average and standard deviation of the performance over all the models.

Additional experiments were also conducted on another dialogue domain. We use the SACTI dialogue management dataset (Williams & Young, 2004), which consists of a set of conversations between a user and a help wizard over the telephone. Both the user and wizard are human operators. The audio is recorded and transcribed. The speech input is processed through a speech recognizer.

The task domain is one in which the user can request information concerning nine different topics, including restaurants, movies, bus schedules, hotel locations, etc. The user can also request information concerning up to two topics concurrently. For example, he may be asking for the directions to a hotel, thus requesting map information and hotel information. Overall, twenty-five instances of single or pairs of topics appear in the dataset. These constitute the state space for the dialogue management problem. The set of words emitted by the user define the natural observation set. In total, 448 unique words occur in the data set. We assume the user can continue discussing the same topic or change topics after every word. A portion of the data set covering 25 conversations (each consisting of 10-30 utterances) was annotated by hand to have accurate state labellings. The POMDP parameters (topic-to-topic transition probabilities, word emission probabilities) are estimated from the annotated dataset. The reward function was defined as positive when the correct action is taken (i.e. relevant information is displayed). A feature space reduction was applied, resulting in a model with 25 states, 25 actions, and 25 observations (Atrash & Pineau, 2006).

The results using the Bayesian reinforcement framework on the SACTI domain are found in figure 7-5. Parameter tying is not used in this case, although many

observation probabilities are zero. This results in learning approximately 100 observation parameters and 2 reward parameters. A prior is used by determining rough counts based on the ground truth model. As with earlier domains, as more data is collected, the models are able to select the correct action more often. However, because the domain is more complex, the rate of improvement is slower.

7.4 Related Work

Early work (Chrisman, 1992) on the application of Baum-Welch to general POMDPs used the same mechanisms of expected counts to re-estimate the parameters iteratively after data is collected. This work considers the idea of a predictive model and also presented mechanisms to adjust the number of underlying states as needed to maximize the predictive power of the model. Baum-Welch has been used to train POMDPs for robot navigation and localization for the Xavier robot (Koenig & Simmons, 1996b) (Koenig & Simmons, 1996a). Here, POMDPs were used to represent the physical layout of a building. The transition and observation parameters were learned based on traces of the robot moving through the environment using sliding window of history over which the parameters were refined at intervals. The topology of the environment was learned by maintaining several potential distances of each hallway over which the model is learned. Baum Welch was again used to learn navigation POMDPs (Shatkey & Kaelbling, 1997), this time with less constraints on than topology. Both of these works took advantage of the domain to aid in learning. Other research (Nikovski & Nourbakhsh, 2000) addressed the use of topology to aid with learning. In this work, a very large initial model is generated.

As traces are run through the model, heuristics are used to determine if two states should be merged.

7.5 Future Work

The experiments in the previous two chapters focus on the reward and observation parameters, however many domains also require the transition parameters to be learned as well. While the learning of these parameters was not directly addressed in this thesis, the mechanism for learning the observation parameters is directly applicable. The transition parameters can be represented with a Dirichlet distribution, and a similar update to the distribution applied given the information from the oracle.

7.6 Conclusion

In this chapter, we extend the application of the Bayesian reinforcement learning framework to learning observation probabilities. We demonstrate learning over the observation parameters, as well as the observation and reward parameters together. As with the previous chapter, these results demonstrate the framework is able to improve performance over time given information from the oracle when the observation and reward parameters are initially unknown.

We compare the performance of the framework against Baum-Welch, the standard algorithm for learning observation parameters in partially observable domains. While Baum-Welch performed better for learning only observations, we demonstrate that the BRL algorithm, which is able to consider the policy while learning, is robust to noise in the model, particularly in parameters which cannot be altered. We show that in the dialogue domain, by considering the policy, the Bayesian reinforcement

learning approach is able to compensate for the error in the model by adjusting the other parameters. These results suggest there are lessons to be learned from the Baum-Welch approach. For example, while the BRL method updates based on the immediate experiences, Baum-Welch considers a sequence of data, including a history and future from the current time, to perform an update to the parameter.

CHAPTER 8

Conclusion

As robots and autonomous agents become more prevalent in real-world domains, they face increasingly difficult and complex decision-making tasks. These agents must account for uncertainty in the environment due to noisy sensors and actuators while performing robust decision-making. This is particularly true in domains where agents are working alongside humans as the agent must now handle the complexity of human behavior.

The Smartwheeler autonomous wheelchair project presented in this thesis represents a typical domain which must face these issues. In addition to traditional robotics problems, such as navigation, the autonomous wheelchair must handle the uncertainty for the speech recognition interface. The wheelchair must determine the user's intent despite the prevalence of errors by the recognition system, as well as the ambiguity that arises from human language.

This thesis addresses these issues by presenting a framework for learning parameters for Partially Observable Markov Decision Processes. POMDPs have the advantage of capturing both noise from actuation, as well as from sensing, but their use requires accurate models of the environment. The aim of the thesis is to extend the application of POMDPs by providing a mechanism for automatic model learning. By adapting models through experience, agents can deal with rough or incorrect models, as well as environments that change over time.

To address learning, the framework presented draws upon a Bayesian reinforcement learning representation. Additionally, this framework takes advantage of domains where an agent works alongside a human operator to help guide the learning through the use of policy information. To support this framework, experimental results are presented using data collected from users studies on the Smartwheeler autonomous wheelchair.

8.1 Framework

The framework presented in this thesis uses Bayesian reinforcement learning to determine the parameters of a POMDP. While traditional reinforcement learning methods use point estimates, Bayesian reinforcement learning represents the parameters as distributions. This allows the agent to reason about the uncertainty over the learned parameters themselves. This thesis addresses issues which arise by using such a representation, such as the choice of distributions, selection of actions, and updating the distributions given experience. Additionally, while earlier work focuses on fully observable domains, this thesis emphasizes learning in partially observable domains, which poses further issues as the underlying state must be inferred, instead of observed directly.

One of the innovations of this framework is the use of the policy oracle to help guide the search through parameter space. In this case, by using the policy provided by an oracle to help select representative sampled models, the hyper-parameters of the distributions can be moved accordingly towards the correct model. While supplying optimal actions provides less information to the framework than state information, it may be easier for the user to provide in many situations. Requiring

alternate information such as state information may be difficult in domains such as dialogue management.

While much of the earlier work on POMDP learning has focused on the observation and transition parameters, a significant advantage of the framework in this thesis is the ability to learn the reward parameters. The reward parameters are often the most difficult to specify especially when not physically grounded, unlike transitions and observations. The reward values allow a model to capture a relationship between the actions themselves, and as such have a significant impact on the resulting policy. This framework is able to determine reward parameters based on the policy provided by the oracle. Furthermore, while previous work has focused on learning only reward or observation parameters, this framework is able to learn both sets of parameters in unison.

8.2 Experimental Results

A major contribution of this thesis is the validation of the framework using real-world data collected from user studies of the Smartwheeler autonomous wheelchair. Earlier work typically validates algorithms on artificial or very simple MDP and POMDP domains. The results presented in this thesis focus on data collected from authentic users performing a user study of the wheelchair in a rehabilitation center environment. The control interface for the robot wheelchair provides a challenging domain which requires this type of framework due to the combination of misrecognition and ambiguity of the speech recognition and need for the system to infer the user's intentions. This rich and complex domain requires extensive expert knowledge, as well as trial and error, to implement manually, and would vary on a user-by-user

basis. Providing a framework capable of adapting the provided rough models to each user is necessary for such a system to be viable.

The results presented in this thesis demonstrate an improvement in the performance of the framework as more data is collected from the oracle. This is demonstrated when learning the reward and observation parameters, separately and concurrently. The results also demonstrate the various aspects of the framework, including the selection of models for performing the update, as well as the impact of the sampling and the update to the models.

Finally, these results are compared against an instance-based learning approach as well as the Baum-Welch algorithm. The instance-based learning approach was unable to capture the temporal aspects of the domain required to perform decision-making. This method does not return a model, thus any policy learned would be difficult to generalize beyond the scenario in which it was learned. Although models trained using Baum-Welch perform better initially when only learning the observation parameters, our framework is able to approach the performance with additional data. Baum-Welch, however, lacks any mechanism to learn the reward parameter.

8.3 Limitations and Future Work

An issue which has not been addressed up to this point is the number of queries to the oracle. The previous results have shown the effects of training on the system by querying the oracle at every step, yet as shown in the figures, the method requires a very large number of queries. This is not unusual for reinforcement learning, where millions of experience steps are sometimes required to learn even simple tasks (Sutton

& Barto, 1998). However this is clearly problematic for our target domain, given that we presume the oracle to be a human user or operator.

Previous work (Doshi et al., 2008; Jaulmes et al., 2005) has used heuristics for selecting when to perform a query. These heuristics include examining the entropy of the current belief state, the disagreement of action selection over the models, as well as the Bayes risk associated with action. These heuristics are directly applicable to the framework presented in this thesis and can be used to limit the queries issued.

An additional approach, adapted from standard (non-Bayesian) reinforcement learning, is to remember the previous oracle queries and refer to them before querying the oracle directly. This requires that the optimal actions for previously visited belief states be stored in memory. If a similar belief state is retrieved from memory, the oracle action recommended for that memorized belief state is applied to the new belief state. By introducing this supervised learning element into the system, we would be able to take advantage of the increased speed of supervised learning, while at the same time maintaining the generalizability of the underlying model. Examples of using supervised learning to train policies include modelling a policy using a neural network driving an autonomous vehicle (Pomerleau, 1995) and using a nearest neighbor method to learning options (Păduraru & Bulitko, 2005). This is also similar to work on batch reinforcement learning (Kalyanakrishnan & Stone, 2007; Lin, 1992; Ernst et al., 2005) where experiences encountered by an agent are stored and used multiple times in the learning process. This would help speed up learning, as the agent would normally have to re-encounter those experiences repeatedly under the original reinforcement learning methods.

The validation in this work centered around two versions of a dialogue control system, a *data-based* model using belief states collected from experiments, and a *generative-based* model with belief states generated using the POMDP itself. Although both models utilized the same underlying POMDP, the system trained using data collected from users was able to learn the parameters more quickly. The significant difference in the training of these two models comes from the number of different belief states encountered, with the *data-based* approach encountering significantly more belief states. This suggests there may be a potential characterization of the data that would provide insight about the quality of the models that can be learned.

The use of an oracle and the selection of heuristics is closely tied to the human-robot interaction component to the system. Ideally, any improvements to the oracle would be validated with user studies to understand the proper heuristics to apply. These studies would not only examine the performance of the system, but also survey the user operating the system. This would help further explore the most comfortable way for a user to provide feedback. For example, there may be alternative methods to obtain information from the user. The oracle in this thesis assumes an action is specified by the oracle. An alternative approach may be to select and execute an action unless the user corrects the oracle. There may be many situations where many actions are acceptable to the user or the user has no preference. By taking advantage of this “negative” oracle, the user interaction is limited only to the critical situations important to the user.

Finally, in addition to limiting the queries, it is important to take advantage of the data which is collected. One improvement relies on batch processing to reduce the required data. When deployed to a real-time environment, the framework could store experiences in a manner similar to the experience re-use mechanism or the iterations used in Baum-Welch, and perform iterations of the framework with the stored data. This would help take further advantage of the data collected from the user. One potential direction is to consider an approach in the direction of an expectation-maximization update similar to Baum-Welch. Future work may address combining online learning of POMDP parameters through Baum-Welch (Koenig & Simmons, 1996b) with a Bayesian framework to allow for more robust decision-making and the ability to learn reward parameters.

8.4 Summary

POMDPs are powerful probabilistic models which are well suited for many problems faced by autonomous agents. However, they require a method for learning parameters automatically to perform in more complex real-world domains. Through learning, POMDPs can be applied to novel domains which may otherwise be too complex to model manually

This thesis presents a novel method for learning those parameters using Bayesian reinforcement learning in situations with an agent working alongside a human. By leveraging policy information provided by a user, the framework can adapt the model parameters to generate policies which replicate that of the user. Requiring only policy information allows the user to guide the system by demonstration. This capability

allows agents to operate in domains which would otherwise be too complex or difficult to model, and facilitates humans and machines working in coordination.

Appendix A - Distributions

The Beta distribution is a continuous probability distribution over the interval $[0, 1]$. This distribution is defined by two hyper-parameters, typically called α and β . The Beta distribution acts as the conjugate prior for binomial distributions, effectively acting as the prior over the relationship of two random variables. The probability distribution function is defined as:

$$pdf(x, \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$$

where B is the beta function:

$$B(\alpha, \beta) = \frac{(\alpha-1)!(\beta-1)!}{(\alpha+\beta-1)!}$$

and $\alpha, \beta > 0$. The expected value and variance are defined as:

$$E(X) = \frac{\alpha}{\alpha + \beta}$$

$$Var(X) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

The Dirichlet distribution is a continuous multivariate probability distribution parameterized by a set of hyper-parameters $\alpha_0, \alpha_1, \dots, \alpha_N$ and is commonly used as a prior distribution in Bayesian reasoning. This distribution is the conjugate prior of the multivariate distribution, acting as the prior over the relationship of several random variables. The Beta distribution is a special case of the Dirichlet distribution for only two parameters. The probability distribution function for a Dirichlet

distribution is defined as:

$$p(x_0, x_1, \dots, x_N; \alpha_0, \alpha_1, \dots, \alpha_N) = \frac{\Gamma(\sum_{i=0}^N \alpha_i)}{\prod_{i=0}^N \Gamma(\alpha_i)} \prod_{i=0}^N x_i^{\alpha_i - 1}$$

where Γ is the Gamma function, a generalized factorial function for real numbers:

$$\Gamma(n) = (n - 1)!$$

and $\alpha_i > 0$. The expected value and variance are defined as:

$$E[X_j] = \frac{\alpha_j}{\sum_{i=0}^N \alpha_i}$$
$$Var(X_j) = \frac{E[X_j](1 - E[X_j])}{\sum_{i=0}^N \alpha_i + 1}$$

Appendix B - Interaction Manager POMDP

The model used to control the dialogue in the interaction studies was required to handle 28 potential user intentions:

- align to wall
- align to wall left
- align to wall right
- approach door
- drive backward
- drive backward a little
- drive forward
- drive forward a little
- drive forward X meters
- go through door
- park to the left
- park to the right
- set mode to indoor
- set mode to outdoor
- set speed to fast
- set speed to medium
- set speed to slow
- stop
- turn around
- turn left
- turn left a little
- turn left X degrees
- turn right
- turn right a little
- turn right X degrees
- veer left
- veer right
- move forward along the slope

To represent these commands, a partially observable Markov decision process was defined as follows:

- **States, S :** The interaction POMDP has 28 states, one state, s_0, s_1, \dots, s_{27} for each potential user intention.
- **Actions, A :** The model includes one action, a_0, a_1, \dots, a_{27} , for each potential user intention, plus 2 specific queries, a_{28}, a_{29} , and one general query, a_{30} . Each specific query is related to one of the two classes of actions, in this case one query for the movement actions, and one query for the configuration actions.
- **Observations, Z :** Observations are received from the semantic parser and consist of assignments to the semantic slots

$$z = (CommandValue = v_1, CommandType = v_2, DirectionValue = v_3, \dots) \quad (8.1)$$

If there is no valid parse of the speech, the observation consists of a bag-of-words, $z = (word_1 = v_1, word_2 = v_2, word_3 = v_3, \dots)$.

- **Rewards:** Since each state is associated with a user intention, there is one “correct” action to take in each state. Each user intention is also associated with one of the two specific queries, depending on the type of action. The rewards, $R(s, a)$ were found by hand through repeated trials:

$$R(s, a) = \begin{cases} 0 & \text{if } a \in \{a_0, a_1, \dots, a_{27}\} \text{ and } a \text{ is not correct action for } s \\ 1 & \text{if } a \in \{a_0, a_1, \dots, a_{27}\} \text{ and } a \text{ is correct action for } s \\ 0.3 & \text{if } a \in \{a_{28}, a_{29}\} \text{ and } a \text{ is not correct specific query for } s \\ 0.7 & \text{if } a \in \{a_{28}, a_{29}\} \text{ and } a \text{ is correct specific query for } s \\ 0.5 & \text{if } a = a_{30} \end{cases} \quad (8.2)$$

- **Transition Distribution, T :** When the agent executes a non-query action, it is assumed the user has a new intention. The model transitions from state to state with uniform probability. When a query is issued, it is assumed the

intention of the user remains the same, thus the state does not change. More formally:

$$T(s_i, a_j, s_k) = \begin{cases} 1/28 & \text{if } a \in \{a_0, a_1, \dots, a_{27}\} \\ 0 & \text{if } a \in \{a_{28}, a_{29}, a_{30}\}, s_i \neq s_k \\ 1 & \text{if } a \in \{a_{28}, a_{29}, a_{30}\}, s_i = s_k \end{cases} \quad (8.3)$$

- **Observation Distribution:** The assignment of values to semantic slots is assumed to be independent. Thus:

$$P(s_i|v_1, v_2, v_3, \dots) = \frac{P(v_1, v_2, v_3, \dots|s_i)P(s_i)}{P(v_1, v_2, v_3, \dots)} \quad (8.4)$$

$$= \frac{P(v_1|s_i)P(v_2|s_i)P(v_3|s_i)\dots P(s_i)}{P(v_1)P(v_2)P(v_3)\dots} \quad (8.5)$$

$P(v_j), P(s_i)$, and $P(v_j|s_i)$ are determined through labeled data collected via experimentation. Phrases spoken by users were labeled based on the users true intention which were then used to calculate the necessary probabilities.

Bibliography

- Abbeel, P., & Ng, A. Y. (2004). Apprenticeship learning via inverse reinforcement learning. *International Conference on Machine learning*.
- Amori, R. D. (1992). Vocomotion - an intelligent voice-control system for powered wheelchairs. *Rehabilitation Engineering and Assistive Technology Society of North America* (pp. 421–423).
- Anderson, B., & Moore, A. (2005). Active learning for hidden Markov models: Objective functions and algorithms. *International Conference on Machine Learning* (pp. 9–16).
- Armstrong-Crews, N., & Veloso, M. (2007). Oracular partially observable Markov decision processes: a very special case. *International Conference on Robotics and Automation*. Rome, Italy.
- Atkeson, C., & Schaal, S. (1997). Robot learning from demonstration. *International Conference on Machine Learning* (pp. 12–20).
- Atrash, A., Kaplow, R., Villemure, J., West, R., Yamani, H., & Pineau, J. (2009). Development and validation of a robust speech interface for improved human-robot interaction. *International Journal of Social Robotics*, 1, 345–356.
- Atrash, A., & Pineau, J. (2006). Efficient planning and tracking in POMDPs with large observation spaces. *AAAI Workshop on Empirical and Statistical Approaches for Spoken Dialogue Systems* (pp. 7–12).

- Atrash, A., & Pineau, J. (2009). A Bayesian reinforcement learning approach for customizing human-robot interfaces. *International Conference on Intelligent User Interfaces* (pp. 355–360).
- Bagnell, J., & Schneider, J. (2001). Autonomous helicopter control using reinforcement learning policy search methods. *International Conference on Robotics and Automation* (pp. 1615–1620).
- Balentine, B., & Morgan, D. (2001). *How to build a speech recognition application: Second edition: A style guide for telephony dialogues*. Enterprise Integration Group.
- Barto, A., Sutton, R., & Anderson, C. (1983). Neuron-like adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics* (pp. 834–846).
- Bellman, R. (1957). *Dynamic programming*. Dover Publications.
- Bertsekas, D., & Tsitsiklis, J. (1996). *Neurodynamic programming*. Belmont, MA: Athena-Scientific.
- Cassandra, A. (1998). A survey of POMDP applications. *AAAI Fall Symposium on POMDPs*.
- Cassandra, A., Littman, M. L., & Zhang, N. (1997). Incremental pruning: A simple, fast, exact method for partially observable Markov decision processes. *Uncertainty in Artificial Intelligence* (pp. 54–61).
- Chalkiadakis, G., & Boutilier, C. (2003). Coordination in multiagent reinforcement learning: a Bayesian approach. *Autonomous Agents and Multiagent Systems* (pp. 709–716).

- Cheng, H.-T. (1988). *Algorithms for partially observable Markov decision processes*.
Doctoral dissertation, University of British Columbia.
- Choi, J., & Kim, K. E. (2009). Inverse reinforcement learning in partially observable environments. *International Joint conference on Artificial intelligence* (pp. 1028–1033). Pasadena, California, USA.
- Chrisman, L. (1992). Reinforcement learning with perceptual aliasing: The perceptual distinctions approach. *National Conference on Artificial Intelligence* (pp. 183–188).
- Cohen, M., Giangola, J., & Balogh, J. (2004). *Voice user interface design*. Addison-Wesley Professional.
- Collett, T. H. J., MacDonald, B. A., & Gerkey, B. P. (2005). Player 2.0: Toward a practical robot programming framework. *Australasian Conference on Robotics and Automation*.
- Davis, A., Moore, M., & Storey, V. (2003). Context-aware communication for severely disabled users. *Conference on Universal Usability* (pp. 106–111).
- Dearden, R., Friedman, N., & Andre, D. (1999). Model based Bayesian exploration. *Uncertainty in Artificial Intelligence* (pp. 150–159).
- Dearden, R., Friedman, N., & Russell, S. (1998). Bayesian Q-learning. *American Association of Artificial Intelligence* (pp. 761–768).
- Demeester, E., Huntemann, A., Vanhooydonck, D., Vanacker, G., Degeest, A., Brusel, H. V., & Nuttin, M. (2006). Bayesian estimation of wheelchair driver intents: Modeling intents as geometric paths tracked by the driver. *International Conference on Intelligent Robots and Systems* (pp. 5775 – 5780).

- Dempster, A., Laird, N., & Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of Royal Statistical Society*, 39, 1–38.
- Doran, C., Aberdeen, J., Damianos, L., & Hirschman, L. (2001). Comparing several aspects of human-computer and human-human dialogues. *SIGdial Workshop on Discourse and Dialogue* (pp. 1–10).
- Doshi, F., Pineau, J., & Roy, N. (2008). Reinforcement learning with limited reinforcement: Using Bayes risk for active learning in POMDPs. *International Symposium on Artificial Intelligence and Mathematics*.
- Doshi, F., & Roy, N. (2007). Efficient model learning for dialog management. *International Conference on Human-Robot Interaction* (pp. 65–72).
- Duff, M. (2002). *Optimal learning: Computational procedures for bayes-adaptive markov decision processes*. Doctoral dissertation, University of Massachusetts Amherst. Director-Andrew Barto.
- Ernst, D., Geurts, P., & Wehenkel, L. (2005). Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6, 503–556.
- Even-Dar, E., Kakade, S., & Mansour, Y. (2005). Reinforcement learning in POMDPs without resets. *International Joint Conference on Artificial Intelligence* (pp. 690–695). Edinburgh, Scotland.
- Evermann, G., & Woodland, P. (2000). Posterior probability decoding, confidence estimation and system combination. *Speech Transcriptions Workshop*.
- Fehr, L., Langbein, W., & Skaar, S. (2000). Adequacy of power wheelchair control interfaces for persons with severe disabilities: A clinical survey. *Journal of Rehabilitation Research & Development*, 37(3), 353–360.

- Felzer, T., & Freisleben, B. (2002). HaWCoS:the 'hands-free' wheelchair control system. *International ACM conference on Assistive technologies* (pp. 127–134).
- Gabsdil, M., & Lemon, O. (2004). Combining acoustic and pragmatic features to predict recognition performance in spoken dialogue systems. *Association for Computational Linguistics* (pp. 344–351).
- Goddeau, D., & Pineau, J. (2000). Fast reinforcement learning of dialog strategies. *International Conference on Acoustics, Speech, and Signal Processing*.
- Guestrin, C., Koller, D., & Parr, R. (2001). Multiagent planning with factored mdps. *Neural Information Processing Systems* (pp. 1523–1530).
- Gulati, S., & Kuipers, B. (2008). High performance control for graceful motion of an intelligent wheelchair. *International Conference on Robotics and Automation* (pp. 3932 – 3938).
- Hauskrecht, M. (2000). Value-function approximations for partially observable markov decision processes. *Journal of Artificial Intelligence Research*, 13, 33–94.
- Henderson, J., Lemon, O., & Georgila, K. (2005). Hybrid reinforcement/supervised learning for dialogue policies from communicator data. *Workshop on Knowledge and Reasoning in Practical Dialogue Systems*.
- Hoey, J., Poupart, P., Boutilier, C., & Mihailidis, A. (2005). POMDP models for assistive technology. *AAAI Fall Symposium on Caring Machines: AI in Eldercare*.
- Izadi, M., & Precup, D. (2003). A planning algorithm for predictive state representations. *International Joint Conference on Artificial Intelligence* (pp. 1520–1521).
- Jaakkola, T., Singh, S., & Jordan, M. (1995). Reinforcement learning algorithm for partially observable Markov decision problems. *Advances in Neural Information*

- Processing Systems* (pp. 345–352).
- Jaulmes, R., Pineau, J., & Precup, D. (2005). Active learning in partially observable Markov decision processes. *European Conference on Machine Learning*.
- Kaelbling, L., Littman, M., & Cassandra, A. (1998). Planning and acting in partially observable stochastic domains. *Artificial Intelligence* (pp. 99–134).
- Kaelbling, L., Littman, M., & Moore, A. (1996). Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 237–285.
- Kalyanakrishnan, S., & Stone, P. (2007). Batch reinforcement learning in a complex domain. *International Joint Conference on Autonomous Agents and Multiagent Systems* (pp. 1–8).
- Katevas, N. I., Sgouros, N. M., Tzafestas, S. G., Papakonstantinou, G., Beattie, P., Bishop, J. M., Tsanakas, P., & Koutsouris, D. (1997). The autonomous mobile robot scenario: A sensor-aided intelligent navigation system for powered wheelchairs. *Robotics and Automation Magazine*, 60–70.
- Kemp, T., & Schaaf, T. (1997). Estimating confidence using word lattices. *EuroSpeech* (pp. 827–830).
- Kesrns, M., Mansour, Y., & Ng, A. (2001). A sparse sampling algorithm for near-optimal planning in large Markov decision process. *Machine Learning*, 49, 193–208.
- Kirby, R. L. (2007). Wheelchair skills program (WSP), version 4.1, wheelchair skills test (WST) manual. http://www.wheelchairskillsprogram.ca/eng/4.1/WST_Manual_Version_4.1.pdf.

- Kirby, R. L., Dupuis, D. J., & et. al., A. H. M. (2004). The wheelchair skills test (version 2.4): Measurement properties. *Archives of Physical Medicine and Rehabilitation*, 85, 794–804.
- Koenig, S., & Simmons, R. (1996a). Passive distance learning for robot navigation. *International Conference on Machine Learning* (pp. 266 – 274).
- Koenig, S., & Simmons, R. (1996b). Unsupervised learning of probabilistic models for robot navigation. *International Conference on Robotics and Automation*.
- Kurniawati, H., Du, Y., Hsu, D., & Lee, W. (2009). Motion planning under uncertainty for robotic tasks with long time horizons. *International Symposium on Robotics Research*.
- Kurniawati, H., Hsu, D., & Lee, W. S. (2008). SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. *Robotics: Science and Systems*.
- Levin, E., & Pieraccini, R. (1997). A stochastic model of computer-human interaction for learning dialogue strategies. *EUROSPEECH* (pp. 1883–1886).
- Levin, E., Pieraccini, R., & Eckert, W. (1998). Using Markov decision process for learning dialogue strategies. *International Conference on Acoustics, Speech, and Signal Processing* (pp. 201–204).
- Levine, S. P., Bell, D. A., Jaros, L. A., Simpson, R. C., & Koren, Y. (1999). The NavChair assistive wheelchair navigation system. *IEEE Transactions of Rehabilitation Engineering* (pp. 443–451).
- Lin, L.-J. (1992). Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8, 293–321.

- Litman, D., & Pan, S. (2000). Predicting and adapting to poor speech recognition in a spoken dialogue system. *National Conference on Artificial Intelligence* (pp. 722–728).
- Littman, M., Sutton, R., & Singh, S. (2001). Predictive representations of state. *Neural Information Processing Systems* (pp. 1555–1561).
- Littman, M. L. (1996). *Algorithms for sequential decision making*. Doctoral dissertation, Brown University.
- Loch, J., & Singh, S. (1998). Using eligibility traces to find the best memoryless policy in partially observable Markov decision processes. *International Conference on Machine Learning* (pp. 323–331).
- Lockerd, A., & Breazeal, C. (2004). Tutelage and socially guided robot learning. *Intelligent Robots and Systems* (pp. 3475–3480).
- Lovejoy, W. S. (1991). A survey of algorithmic methods for partially observable Markov decision processes. *Annals of Operations Research*, 28, 47–66.
- MacKay, D. (1992). Information-based objective functions for active data selection. *Neural Computation*, 4, 590–604.
- Mandel, C., Huebner, K., & Vierhuff, T. (2005). Towards an autonomous wheelchair: Cognitive aspects in service robotics. *Towards Autonomous Robotic Systems* (pp. 165–172).
- Marc, E., Swerts, M., Theune, M., & Weegels, M. (1999). Problem spotting in human-machine interaction. *EUROSPEECH* (pp. 1423–1426).
- McCallum, A. (1994). Instance-based state identification for reinforcement learning. *Advances in Neural Information Processing Systems* (pp. 377–384).

- McCallum, A. (1995). Instance-based utile distinctions for reinforcement learning. *International Machine Learning Conference* (pp. 387–395).
- Meuleau, N., Hauskrecht, M., Kim, K.-E., Peshkin, L., Kaelbling, L., Dean, T., & Boutilier, C. (1998). Solving very large weakly coupled markov decision processes. *AAAI* (pp. 165–172).
- Mihailidis, A., Elinas, P., Boger, J., & Hoey, J. (2005). An intelligent powered wheelchair to enable mobility of cognitively impaired older adults: An anti-collision system. *IEEE Transactions on Neural Systems & Rehabilitation Engineering* (pp. 136–143).
- Mitchell, T. M. (1997). *Machine learning*. New York: McGraw-Hill.
- Montemerlo, M., Roy, N., & Thrun, S. (2003). Perspectives on standarization in mobile robot programming: The Carnegie Mellon Navigation (CARMEN) Toolkit. *International Conference on Robotics and Systems* (pp. 2436–2441).
- Moore, A., & Atkeson, C. (1993). Prioritized sweeping: Reinforcement learning with less data and less real time. *Machine Learning*, 13, 103–130.
- Ng, A., & Russell, S. (2000). Algorithms for inverse reinforcement learning. *International Conference on Machine Learning* (pp. 663–670).
- Nikovski, D., & Nourbakhsh, I. (2000). Learning probabilistic models for decision-theoretic navigation of mobile robots. *International Conference on Machine Learning* (pp. 671–678).
- Nuance (2009). Dragon naturally speaking 10. <http://www.nuance.com>.
- Nuttin, M., Demeester, E., Vanhooydonck, D., & Brussel, H. V. (2001). Shared autonomy for wheel chair control: Attempts to assess the user’s autonomy. *Autonome*

- Mobile Systeme* (pp. 127–133).
- Păduraru, C., & Bulitko, V. (2005). Supervised Learning of Options: A Pilot Study. *IJCAI Workshop on Planning and Learning in A Priori Unknown or Dynamic Domains* (pp. 155 – 158).
- Paek, T., & Horvitz, E. (2000). Conversation as action under uncertainty. *Uncertainty in Artificial Intelligence* (pp. 455–464).
- Paek, T., Horvitz, E., & Ringer, E. (2000). Continuous listening for unconstrained spoken dialog. *International Conference on Spoken Language Processing*.
- Pao, C., Schmid, ., & Glass, J. (1998). Confidence scoring for speech understanding systems. *International Conference on Speech and Language Processing*.
- Paquet, S., Tobin, L., & Chaib-draa, B. (2005). An online POMDP algorithm for complex multiagent environments. *International joint conference on Autonomous agents and multiagent systems* (pp. 970–977).
- Philips, J., del R. Millan, J., Vanacker, G., Lew, E., Galan, F., Ferrez, P., & Van Brussel, H. and Nuttin, M. (2007). Adaptive shared control of a brain-actuated simulated wheelchair. *International Conference on Rehabilitation Robotics* (pp. 408 – 414).
- Pietquin, O., & Renal, S. (2002). ASR system modeling for automatic evaluation and optimization of dialogue systems. *International Conference on Acoustics, Speech, and Signal Processing* (pp. 46–49).
- Pineau, J., & et. al. (2007). Smartwheeler: A robotic wheelchair test-bed for investigating new models of human-robot interaction. *AAAI Spring Symposium on Multidisciplinary Collaboration for Socially Assistive Robotics* (pp. 59–64).

- Pineau, J., Gordon, G., & Thrun, S. (2003). Point-based value iteration: An anytime algorithm for POMDPs. *International Joint Conference on Artificial Intelligence* (pp. 1025–1032).
- Pires, G., & Nunes, U. (2002). A wheelchair steered through voice commands and assisted by a reactive fuzzy-logic controller. *Journal of Intelligent Robotic Systems*, *34*(3), 301–314.
- Pomerleau, D. (1995). Neural network vision for robot driving. In M. Arbib (Ed.), *The Handbook of Brain Theory and Neural Networks*.
- Poupart, P., & Vlassis, N. (2008). Model-based bayesian reinforcement learning in partially observable domain. *International Symposium on Artificial Intelligence and Mathematics*.
- Poupart, P., Vlassis, N., Hoey, J., & Regan, K. (2006). An analytic solution to discrete bayesian reinforcement learning. *International Conference on Machine Learning* (pp. 697–704). Pittsburgh, Pennsylvania.
- Price, B., & Boutilier, C. (1999). Implicit imitation in multiagent reinforcement learning. *International Conference on Machine Learning* (pp. 325–334).
- Pruski, A., Ennaji, M., & Morere, Y. (2002). VAHM: a user adapted intelligent wheelchair. *International Conference on Control Applications* (pp. 784–789).
- Puterman, M. (1994). *Markov decision processes: Discrete stochastic dynamic programming*. New York, NY: John Wiley and Sons.
- Rabiner, L. R. (1990). A tutorial on hidden Markov models and selected applications in speech recognition. *Readings in Speech Recognition* (pp. 267–296).

- Ramachandran, D., & Amir, E. (2007). Bayesian inverse reinforcement learning. *International Joint Conference on Artificial Intelligence* (pp. 2586–2591).
- Ratliff, N., Bagnell, J., & Zinkevich, M. (2006). Maximum margin planning. *International Conference on Machine Learning* (pp. 729–736).
- Ross, S., Chaib-draa, B., & Pineau, J. (2007). Bayes-adaptive POMDPs. *Neural Information Processing Systems*.
- Roy, N., Pineau, J., & Thrun, S. (2000). Spoken dialogue management using probabilistic reasoning. *Association for Computational Linguistics*.
- Rummery, G., & Niranjan, M. (1994). *On-line q-learning using connectionist systems* (Technical Report).
- Shani, G., & Brafman, R. (2004). Resolving perceptual aliasing in the presence of noisy sensors. *Neural Information Processing Systems* (pp. 1249–1256).
- Shani, G., Brafman, R., & Shimony, S. (2005). Model-based online learning of pomdps. *European Conference on Machine Learning* (pp. 363–364).
- Shani, G., Brafman, R., & Shimony, S. E. (2007). Forward search value iteration for POMDPs. *International Joint conference on Artificial intelligence*.
- Shatkey, H., & Kaelbling, L. (1997). Learning topological maps with weak local odometric information. *International Joint Conference on Artificial Intelligence* (pp. 920–929).
- Simpson, R. (2005). Smart wheelchairs: a literature review. *Journal of Rehabilitation Research & Development*, 42(4), 423–438.
- Simpson, R., Daniel, P., & Baxter, F. (2002). The Hephaestus smart wheelchair system. *IEEE Transactions on neural systems & rehabilitation engineering*, 10(2),

118–122.

- Simpson, R., LoPresti, E., Hayashi, S., Nourbakhsh, I., & Miller, D. (2004). The smart wheelchair component system. *Journal of Rehabilitation Research and Development*, *41(3B)*, 429–442.
- Simpson, R., LoPresti, E. F., & Cooper, R. A. (2008). How many people would benefit from a smart wheelchair? *Journal of Rehabilitation Research & Development*, *45*, 53–72.
- Singh, S., James, M., & Rudary, M. (2004). Predictive state representations: A new theory for modeling dynamical systems. *Uncertainty in Artificial Intelligence* (pp. 512–519).
- Singh, S., Litman, D., Kearns, M., & Walker, M. (2002). Optimizing dialogue management with reinforcement learning: Experiments with the NJFun system. *Journal of Artificial Intelligence Research*, *16*, 105–133.
- Singh, S., Littman, M., Jong, N., Pardoe, D., & Stone, P. (2003). Learning predictive state representations. *International Conference on Machine Learning*.
- Smallwood, R., & Sondik, E. (1973). The optimal control of partially observable markov processes over a finite horizon. *Operations Research*, *21*, 1071–1088.
- Smith, T., & Simmons, R. (2004). Heuristic search value iteration for POMDPs. *Uncertainty in Artificial Intelligence*.
- Sondik, E. (1971). *The optimal control of partially observable Markov decision processes*. Doctoral dissertation, Stanford University.
- Song, L., Boots, B., Siddiqi, S. M., Gordon, G. J., & Smola, A. J. (2010). Hilbert space embeddings of hidden Markov models. *International Conference on Machine*

Learning.

- Spaan, M., & Vlassis, N. (2005). Perseus: Randomized point-based value iteration for POMDPs. *Journal of Artificial Intelligence Research* (pp. 195–220).
- Strens, M. (2000). A Bayesian framework for reinforcement learning. *International Conference on Machine Learning* (pp. 943–950).
- Sutton, R. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1), 9–44.
- Sutton, R. (1990). Integrated architecture for learning, planning, and reacting based on approximating dynamic programming. *International Conference on Machine Learning* (pp. 216–224).
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. MIT Press.
- Szepesvári, C. (2001). Efficient approximate planning in continuous space Markovian decision problems. *AI Communications*, 14, 163–176.
- Szepesvari, C. (2010). *Algorithms for reinforcement learning*. Morgan & Claypool.
- Taha, T., Miro, J., & Dissanayake, G. (2008). POMDP-based long-term user intention prediction for wheelchair navigation. *International Conference on Robotics and Automation* (pp. 3920 – 3925).
- Theocharous, G., Rohanimanesh, K., & Mahadevan, S. (2001). Learning hierarchical partially observable Markov decision process models for robot navigation. *International Conference on Robotics & Automation* (pp. 511–516).

- Tsui, K., & Yanco, H. (2007). Simplifying wheelchair mounted robotic arm control with a visual interface. *AAAI Spring Symposium on Multidisciplinary Collaboration for Socially Assistive Robots* (pp. 97–102).
- Vanhooydonck, D., Demeester, E., Nuttin, M., & Brussel, H. V. (2003). Shared control for intelligent wheelchairs: an implicit estimation of the user intention. *International Workshop on Advances in Service Robotics* (pp. 176–182).
- Wang, T., Lizotte, D., Bowling, M., & Schuurmans, D. (2005). Bayesian sparse sampling for on-line reward optimization. *International Conference on Machine Learning* (pp. 961–968).
- Watkins, C. (1989). *Learning from delayed rewards*. Doctoral dissertation, University of Cambridge.
- White, C. (1991). A survey of solution techniques for the partially observed Markov decision process. *Annals of Operations Research*, 32, 215–230.
- White, M. (2001). Open CCG: The open NLP CCG library. *Software available at <http://openccg.sourceforge.net>*.
- Williams, J. (2008). The best of both worlds: Unifying conventional dialog systems and pomdps. *Interspeech*.
- Williams, J., Poupart, P., & Young, S. (2005a). Factored partially observable markov decision processes for dialogue management. *Workshop on Knowledge and Reasoning in Practical Dialog Systems* (pp. 76–82).
- Williams, J., Poupart, P., & Young, S. (2005b). Partially observable Markov decision processes with continuous observations for dialogue management. *SigDial Workshop on Discourse and Dialogue* (pp. 393–422).

- Williams, J., & Young, S. (2004). Characterizing task-oriented dialog using a simulated ASR channel. *International Conference on Speech and Language Processing*.
- Williams, J., & Young, S. (2005). Scaling up POMDPs for dialog management: The “summary POMDP” method. *Automatic Speech Recognition and Understanding*, 177–182.
- Williams, J., & Young, S. (2006). Scaling POMDPs for dialog management with composite summary point-based value iteration (CSPBVI). *Workshop on Statistical and Empirical Approaches for Spoken Dialogue Systems*.
- Wobbrock, J., Myers, B., Aung, H., & LoPresti, E. (2004). Text entry from power wheelchairs: edgewise for joysticks and touchpads. *Conference on Computers and Accessibility* (pp. 110–117).
- Yanco, H. (2001). Development and testing of a robotic wheelchair system for outdoor navigation. *Conference of the Rehabilitation Engineering and Assistive Technology Society of North America*.
- Zalzal, V., R., Gava, Kelouwani, S., & Cohen, P. (2009). Acropolis: A fast prototyping robotic application. *International Journal of Advanced Robotic Systems*, 6(1), 1–6.
- Zhang, B., Cai, Q., Mao, J., Chang, E., & Guo, B. (2001a). Spoken dialogue management as planning and acting under uncertainty. *Eurospeech*.
- Zhang, B., Cai, Q., Mao, J., & Guo, B. (2001b). Planning and acting under uncertainty: A new model for spoken dialogue systems. *Uncertainty in Artificial Intelligence* (pp. 572–57).