

Efficient Methods for Near-Optimal Sequential Decision Making Under Uncertainty*

Christos Dimitrakakis

April 20, 2010

Abstract

This chapter discusses decision making under uncertainty. More specifically, it offers an overview of efficient Bayesian and distribution-free algorithms for making nearly-optimal sequential decisions under uncertainty about the environment. Due to the uncertainty, such algorithms must not only learn from their interaction with the environment, but also perform as well as possible while learning is taking place.

Contents

1	Introduction	2
1.1	Notation	3
2	Decision making under uncertainty	3
2.1	Utility, randomness and uncertainty	4
2.2	Uncertain outcomes	4
2.3	Bayesian inference	5
2.4	Distribution-free bounds	6
3	Sequential decision making under uncertainty	8
3.1	Stopping problems	8
3.2	Bandit problems	10
3.3	Reinforcement learning and control	10
4	Markov decision processes	11
5	Belief-augmented Markov decision processes	12
5.1	Bayesian inference with a single MDP model class	12
5.2	Constructing and solving BAMDPs	13
5.3	Belief tree expansion	14
5.4	Bounds on the optimal value function	15
5.4.1	Leaf node lower bound	16
5.4.2	Bounds with high probability	17
5.5	Discussion and related work	18

*This is a draft of the chapter appearing in [Dimitrakakis, 2010]. The final publication is available at www.springerlink.com.

6	Partial observability	19
6.1	Belief POMDPs	20
6.2	The belief state	20
6.3	Belief compression	21
7	Conclusion, future directions and open problems	21

1 Introduction

It could be argued that automated decision making is the main application domain of artificial intelligence systems. This includes tasks from selecting moves in a game of chess and choosing the best navigation route in a road network to responding to questions posed by humans, playing a game of poker and exploring other planets of the solar system. While chess playing and navigation both involve accurate descriptions of the problem domain, the latter problems involve uncertainty about both the nature of the environment and its current state. For example, in the poker game, the nature of the other players (i.e. the strategies they use) is not known. Similarly, the state of the game is not known perfectly – i.e. their cards are not known, but it might be possible to make an educated guess.

This chapter shall examine acting under uncertainty in environments with state, wherein a *sequence* of decisions must be made. Each decision made has an effect on the environment, thus changing the environment’s state. One type of uncertainty arises when it is not known how the environment works, i.e. the agent can observe the state of the environment but is not certain what is the effect of each possible action in each state. The decision making agent must therefore explore the environment, but not in a way that is detrimental to its performance. This balancing act is commonly referred to as the exploration-exploitation trade-off. The chapter gives an overview of current methods for achieving nearly optimal online performance in such cases.

Another type of uncertainty arises when the environment’s state can not be observed directly, and can only be inferred. In that case, the state is said to be hidden or partially observable. When both types of uncertainty occur simultaneously, then the problem’s space complexity increases polynomially with time, because we must maintain all the observation history to perform inference. The problem becomes even harder when there are multiple agents acting within the environment. However, we shall not consider this case in this chapter.

The remainder of this chapter is organized as follows. Firstly, we give an introduction to inference and decision making under uncertainty in both the Bayesian and distribution-free framework. Section 3 introduces some sequential decision making problems under uncertainty. These problems are then formalized within the framework of Markov decision processes in Section 4, which can be used to make optimal decisions when the uncertainty is only due to stochasticity, i.e. when the effects of any decisions are random, but arise from a known probability distribution that is conditioned on the agent’s actions and the current state. Section 5 discusses the extension of this framework to when these probability distributions are not known. It is shown that the result is another Markov decision process with an infinite number of states and various methods for approximately solving it are discussed. When the states of the environment are not directly observed, the problem becomes much more complex: this case is examined in Section 6. Finally, Section 7 identifies open problems and possible directions of future research.

1.1 Notation

We shall write $\mathbb{I}\{X\}$ for the indicator function that equals 1 when X is true and 0 otherwise. We consider actions $a \in \mathcal{A}$ and contexts (states, environments or outcomes) $\mu \in \mathcal{M}$. We shall denote a sequence of observations from some set \mathcal{X} as $x^t \triangleq x_1, \dots, x_t$, with $x_k \in \mathcal{X}$.

In general, $\mathbf{P}(X)$ will denote the probability of any event X selected from nature and \mathbf{E} to denote expectations. When observations, outcomes or events are generated via some specific process μ , we shall explicitly denote this by writing $\mathbf{P}(\cdot|\mu)$ for the probability of events. Frequently, we shall use the shorthand $\mu(\cdot)$ to denote probabilities (or densities, when there is no ambiguity) under the process μ . With this scheme, we make no distinction between the name of the process and the distribution it induces. Thus, the notation $\mu(\cdot)$ may imply a marginalisation. For instance, if a process μ defines a probability density $\mu(x, y)$ over observations $x \in \mathcal{X}$, $y \in \mathcal{Y}$, we shall write $\mu(x)$ for the marginal $\int_{\mathcal{Y}} \mu(x, y) dy$. Finally, expectations under the process will be written as $\mathbf{E}_{\mu}(\cdot)$ or equivalently $\mathbf{E}(\cdot|\mu)$. In some cases it will be convenient to employ equality relations of the type $\mu(x_t = x)$, to denote the density at x at time t under process μ .

2 Decision making under uncertainty

Imagine an agent acting within some environment and let us suppose that it has decided upon a fixed plan of action. Predicting the result of any given action within the plan, or the result of the complete plan, may not very easy, since there can be many different sources of uncertainty. Thus, it might be hard to *evaluate* actions and plans and consequently, to find the *optimal* action or plan. This section will focus in the case where only a single decision must be made.

The simplest type of uncertainty arises when events that take place within the world can be stochastic. Those may be (apparently) truly random events, such as which slit will a photon will pass through in the famous two-slit experiment, or events that can be considered random for all practical purposes, such as the outcome of a die roll. A possible decision problem in that setting would be whether to accept or decline a particular bet on the outcome of one die roll: if the odds are favorable, we should accept, otherwise decline.

The second source of uncertainty arises when we do not know exactly how the world works. Consider the problem of predicting the movement of planets given their current positions and velocities. Modelling their orbits as circular, will of course result in different predictions to modelling their orbits as elliptic. In this problem the decision taken involves the selection of the appropriate model.

Estimating which model, or set of models, best corresponds to our observations of the world becomes harder when there is, in addition, some observation stochasticity. In the given example, that would mean that we would not be able to directly observe the planets' positions and thus it would be harder to determine the best model.

The model selection problems that we shall examine in this section involve two separate, well-defined, phases. The first phase involves collecting data, and the second phase involves making a decision about which is the best model. Many statistical inference problems are of this type, such as creating a classifier from categorical data. The usual case is that the observations have already been collected and now form a fixed *dataset*. We then define a set of classification models and the decision making task is to choose one or more classifiers from the given set of models. A lot of recent work on classification algorithms is in fact derived from this type of decision making framework Blumer et al. [1989]; Vapnik [2000]; Vapnik. and Chervonenkis [1971]. A straightforward extension of the problem

to online decision making resulted in the boosting algorithm Freund and Schapire [1997]. The remainder of this section discusses making single decisions under uncertainty in more detail.

2.1 Utility, randomness and uncertainty

When agents (and arguably, humans Savage [1972]) make decisions, they do so on the basis of some preference order among possible outcomes. With perfect information, the rational choice is easy to determine, since the probability of any given outcome given the agent's decision is known. This is a common situation in games of chance. However, in the face of uncertainty, establishing a preference order among actions is no longer trivial.

In order to formalize the problem, we allow the agent to select some action a from a set of \mathcal{A} of possible choices. We furthermore define a set of contexts, states, or environments \mathcal{M} , such that the preferred action may differ depending which is the current context $\mu \in \mathcal{M}$. If the context is perfectly known, then we can simply take the most preferred action in that context.

One way to model this preference is to define a utility function $U : \mathcal{A} \times \mathcal{M} \rightarrow \mathbb{R}$ mapping from the set of possible μ and a to the real numbers.

Definition 2.1 (Utility) *For any context $\mu \in \mathcal{M}$, and actions $a_1, a_2 \in \mathcal{A}$, we shall say that we prefer a_1 to a_2 and write $a_1 \succ a_2$, if and only if $U(a_1, \mu) > U(a_2, \mu)$. Similarly, we write that $a_1 = a_2$ iff $U(a_1, \mu) = U(a_2, \mu)$.*

The transitivity and completeness axioms of utility theory are satisfied by the above definition, since the utility function's range are the real numbers. Thus, if both U and μ are known, then the optimal action must exist and is *by definition* the one that maximizes the utility for the given context μ .

$$a^* = \arg \max_{a \in \mathcal{A}} U(a, \mu).$$

We can usually assign a *subjective preference* for each action a in all μ , thus making the function U well-defined. It possible, however, that we have some uncertainty about μ . We are then obliged to use other formalisms for assigning preferences to actions.

2.2 Uncertain outcomes

We now consider the case where μ is uncertain. This can occur when μ is chosen randomly from some known distribution with density p , as is common in lotteries and random experiments. It may be also chosen by some adversary, which is usual in deterministic games such as chess. In games of chance, such as backgammon, it is some combination of the two. Finally, μ could be neither randomly nor selected by some adversary, but in fact, simply not precisely known: we may only know a set \mathcal{M} which contains μ .

Perhaps the simplest way to assign preferences in the latter case is to select the action with the highest worst-case utility:

Definition 2.2 (Maximin utility) *Our preference $V(a)$ for action a is:*

$$V(a) \triangleq \inf_{\mu \in \mathcal{M}} U(a, \mu). \quad (2.1)$$

This is mostly a useful ordering for the adversarial setting. In the stochastic setting, its main disadvantage is that we may avoid actions which have the highest utility for most high-probability outcomes, apart for some outcomes with near-zero probability, whose utility is small. A natural way to take the probability of outcomes into account, is to use the notion of expected utility:

Definition 2.3 (Expected utility) *Our preference $V(a)$ for action a is the expectation of the utility under the given distribution with density p of possible outcomes:*

$$V(a) \triangleq \mathbf{E}(U|a) = \int_{\mathcal{M}} U(a, \mu) p(\mu) d\mu. \quad (2.2)$$

This is a good method for the case when the distribution from which μ will be chosen is known. Another possible method, which can be viewed as a compromise between expected and maximin utility is to assign preferences to actions based on how likely they are to be close to the best action. For example, we can take the action which has the highest probability of being ε -close to the best possible action, with $\varepsilon > 0$:

Definition 2.4 (Risk-sensitive utility)

$$V(a; \varepsilon) \triangleq \int_{\mathcal{M}} \mathbb{I}\{U(a, \mu) \geq U(a', \mu) - \varepsilon, \forall a' \in \mathcal{A}\} p(\mu) d\mu. \quad (2.3)$$

Thus, an action chosen with the above criterion is guaranteed to be ε -close to the actually best action, with probability $V(a; \varepsilon)$. This criterion could be alternatively formulated as the probability that the action's utility is greater than a fixed threshold θ , rather than being ε -close to the utility of the optimal action. A further modification involves fixing a small probability $\delta > 0$ and then solving for ε , or θ to choose the action which has the lowest regret ε , or the highest guaranteed utility θ , with probability $1 - \delta$. Further discussion of such issues, including some of the above preference relations is given in Friedman and Savage [1948, 1952]; Luce and Raiffa [1957]; Savage [1972].

The above definitions are not strictly confined to the case where μ is random. In *Bayesian*, or *subjectivist* viewpoint of probability, we may also assign probabilities to events which are not random. Those probabilities do not represent possible random outcomes, but *subjective beliefs*. It thus becomes possible to extend the above definitions from uncertainty about random outcomes to uncertainty about the environment.

2.3 Bayesian inference

Consider now that we are acting in one of many possible environments. With knowledge of the true environment and the utility function, it would be trivial, in some sense, to select the utility-maximizing action. However, suppose that we do *not* know which of the many possible environments we are acting in. One possibility is to use the maximin utility rule, but this is usually too pessimistic. An arguably better alternative is to assign a *subjective probability* to each environment, which will represent our belief that it corresponds to reality.¹ It then is possible to use expected utility to select actions.

This is not the main advantage of using subjective probabilities, however. It is rather the fact that we can then use standard probabilistic inference methods to *update our belief* as we acquire more information about the environment. With

¹This is mathematically equivalent to the case where the environment was drawn randomly from a known distribution.

enough data, we can be virtually certain about which is the true environment, and thus confidently take the most advantageous action. When the data is few, a lot of models have a relatively high probability and this uncertainty is reflected in the decision making.

More formally, we define a set of models \mathcal{M} and a prior density ξ_0 defined over its elements $\mu \in \mathcal{M}$. The prior $\xi_0(\mu)$ describes our initial belief that the particular model μ is correct. Sometimes it is unclear how to best choose the prior. The easiest case to handle is when it is known that the environment was randomly selected from a probability distribution over \mathcal{M} with some density ψ : it is then natural (and optimal) to simply set $\xi_0 = \psi$. Another possibility is to use experts to provide information: then the prior can be obtained via formal procedures of prior elicitation Chen et al. [1999]; Dey et al. [1998]. However, when this is not possible either due to the lack of experts or due to the fact that the number of parameters to be chosen is very large, then the priors can be chosen intuitively, according to computational and convenience considerations, or with some automatic procedure: a thorough overview of these topics is given by Berger [2006]; Goldstein [2006]. For now we shall assume that we have somehow chosen a prior ξ_0 .

The procedure for calculating the belief ξ_t at time t is relatively simple: Let x_t denote our observations at time t . For each model μ in our set, we can calculate the posterior probability $\xi_{t+1}(\mu)$ from their prior $\xi_t(\mu)$. This can be done via the definition of joint densities, which is in this form known as Bayes' rule:

$$\xi_{t+1}(\mu) \triangleq \xi_t(\mu|x_t) = \frac{\mu(x_t)\xi_t(\mu)}{\int_{\mathcal{M}} \mu'(x_t)\xi_t(\mu')d\mu'}, \quad (2.4)$$

where we have used the fact that $\xi_t(x_t|\mu) = \mu(x_t)$, since the distribution of observations for a specific model μ is independent of our subjective belief about which models are most likely.

The advantage of using a Bayesian approach to decision making under uncertainty is that our knowledge about the true environment μ at time t is captured via the density $\xi_t(\mu)$, which represents our belief. Thus, we are able to easily utilize any of the action preferences outlined in the previous section by replacing the density p with ξ_t .

As an example, consider the case where we have obtained t observations and must choose the action that appears best. After t observations, we will have reached a belief $\xi_t(\mu)$ and our *subjective* value for each action is simply:

$$V_t(a) \triangleq \mathbf{E}(U|a, \xi_t) = \int_{\mathcal{M}} U(a, \mu)\xi_t(\mu)d\mu. \quad (2.5)$$

At this point, perhaps some motivation is needed to see why this is a good idea. Let μ^* be the true model, i.e. in fact $\mathbf{E}(U|a) = U(a, \mu^*)$ and assume that $\mu^* \in \mathcal{M}$. Then, under relatively lax assumptions², it can be shown (c.f. Savage [1972]) that $\lim_{t \rightarrow \infty} \xi_t(\mu) = \delta(\mu - \mu^*)$, where δ is the Dirac delta function. This implies that the probability measure that represents our belief concentrates³ around μ^* .

There are, of course, a number of problems with this formulation. The first is that we may be unwilling or unable to specify a prior. The second is that the resulting model may be too complicated for practical computations. Finally, although it is relatively straightforward to compute the expected utility, risk-sensitive computations are hard in continuous spaces, since they require calculating the integral of a maximum. In any such situation, distribution-free bounds may be used instead.

²If the true model is not in the set of models, then we may in fact diverge.

³To prove that in more general terms is considerably more difficult, but has been done recently by Zhang Zhang [2006].

2.4 Distribution-free bounds

When we have very little information about the distribution, and we do not wish to utilize “uninformative” or “objective” priors Berger [2006], we can still express the amount of knowledge acquired through observation by the judicious use of distribution-free concentration inequalities. The main idea is that, while we cannot accurately express the possible forms of the underlying distribution, we can always imagine a worst possible case.

Perhaps the most famous such inequality is Markov’s inequality, which holds for any random variable X and all $\varepsilon > 0$:

$$\mathbf{P}(|X| \geq \varepsilon) \leq \frac{\mathbf{E}(|X|)}{\varepsilon}. \quad (2.6)$$

Such inequalities are of greatest use when applied to estimates of unknown parameters. Since our estimates are functions of our observations, which are random variables, the estimates themselves are also random variables. Perhaps the best way to see this is through the example of the following inequality, which applies whenever we wish to estimate the expected value of a bounded random variable by averaging n observations:

Lemma 2.1 (Hoeffding inequality) *If $\hat{x}_n \triangleq \frac{1}{n} \sum_{i=1}^n x_i$, with $x_i \in [b_i, b_i + h_i]$ drawn from some arbitrary distribution f_i and $\bar{x}_n \triangleq \frac{1}{n} \sum_i \mathbf{E}(x_i)$, then, for all $\varepsilon \geq 0$:*

$$\mathbf{P}(|\hat{x}_n - \bar{x}_n| \geq \varepsilon) \leq 2 \exp\left(-\frac{2n^2\varepsilon^2}{\sum_{i=1}^n h_i^2}\right). \quad (2.7)$$

The above inequality is simply interpreted as telling us that the probability of us making a large estimation error decreases exponentially in the number of samples. Unlike the Bayesian viewpoint, where the mean was a random variable, we now consider the mean to be a fixed unknown quantity, and our estimate to be the random variable. So, in practice, such inequalities are useful for obtaining bounds on the performance of algorithms and estimates, rather than expressing uncertainty *per se*.

More generally, such inequalities operate on sets. For any model set \mathcal{M} of interest, we should be able to obtain an appropriate concentration function $\delta(M, n)$ on subsets $M \subset \mathcal{M}$, with $n \in \mathbb{N}$ being the number of observations, such that:

$$\mathbf{P}(\mu \notin M) < \delta(M, n), \quad (2.8)$$

where one normally considers μ to be a fixed unknown quantity and M to be a random estimated quantity. As an example, the right hand side of the Hoeffding inequality can be seen as a specific instance of a concentration function, where $M = \{x : |\hat{x}_n - x| < \varepsilon\}$. A detailed introduction to concentration functions in much more general terms is given by Talagrand [Talagrand, 1996].

An immediate use of such inequalities is to calculate high probability bounds on the utility of actions. Firstly, given the deterministic payoff function $U(a, \mu)$, and a suitable $\delta(M, n)$, let:

$$\theta(M, a) \triangleq \inf_{\mu \in M} U(a, \mu), \quad (2.9)$$

be a lower bound on the payoff of action a in set M . It immediately follows that the payoff $U(a)$ we shall obtain for taking action a will satisfy:

$$\mathbf{P}[U(a, \mu) < \theta(M, a)] \leq \delta(M, n),$$

since the probability of our estimated set not containing μ is bounded by $\delta(M, n)$. Thus, one may fix a set M and then select a maximizing $\theta(a, M)$. This will give a guaranteed performance with probability at least $1 - \delta(M, n)$.

Such inequalities are also useful to bound the expected *regret*. Let the regret of a procedure that has payoff U relative to some other procedure with payoff U^* be $U^* - U$. Now consider that we want take an action a that minimizes the expected regret relative to some maximum value U^* , which can be written as:

$$\mathbf{E}(U^* - U|a) = \mathbf{E}(U^* - U|a, \mu \in M) \mathbf{P}(\mu \in M) + \mathbf{E}(U^* - U|a, \mu \notin M) \mathbf{P}(\mu \notin M), \quad (2.10)$$

where the first term of the sum corresponds to the expected regret that we incur when the model is in the set M , while the right term corresponds to the case when the model is not within M . Each of these terms can be bounded with (2.8) and (2.9): the first term is bounded by $\mathbf{E}(U^* - U|a, \mu \in M) < U^* - \theta(M, a)$ and $\mathbf{P}(\mu \in M) < 1$, while the second term is bounded by $\mathbf{E}(U^* - U|a, \mu \notin M) < U^* - \inf(U)$ and $\mathbf{P}(\mu \notin M) < \delta(M, n)$. Assuming that the infimum exists, the expected regret for any action $a \in \mathcal{A}$ is bounded by:

$$\mathbf{E}(U^* - U|a) \leq (U^* - \theta(M, a)) + (U^* - \inf(U))\delta(M, n). \quad (2.11)$$

Thus, such methods are useful for taking an action which minimizes a bound on the expected regret, that maximizes the probability its utility is greater than some threshold, or finally that maximizes a lower bound on the utility with some at least some probability. However, they cannot be used to select actions that maximize expected utility, simply because the expectation cannot be calculated as we do not have an explicit probability density function over the possible models. We discuss two upper confidence bound based methods for bandit problems in section 3.2 and for the reinforcement learning problem in section 5.5.

3 Sequential decision making under uncertainty

The previous section examined two complementary frameworks for decision making under uncertainty. The task was relatively straightforward, as it involved a fixed period of data collection, followed by a single decision. This is far from an uncommon situation in practice, since a lot of real-world decisions are made in such a way.

In a lot of cases, however, decisions may involve future collection of data. For example, during medical trials, data is continuously collected and assessed. The trial may have to be stopped early if the risk to the patients is deemed too great. A lot of such problems were originally considered in the seminal work of Wald [1947].

This section will present a brief overview of the three main types of sequential decision making problems: Stopping problems, which are the simplest type, bandit problems, which can be viewed as a generalization of stopping problems, and reinforcement learning, which is a general enough framework to encompass most problems in sequential decision making. The reader should also be aware of the links of sequential decision making to classification Freund and Schapire [1997], optimization Auer et al. [2007]; Coquelin and Munos [2007]; Kall and Wallace [1994] and control Agrawal [1995]; Bertsekas [2005, 2001].

3.1 Stopping problems

Let us imagine an experimenter, who needs to make a decision $a \in \mathcal{A}$, where \mathcal{A} is the set of possible decisions. The effect of each different decision will depend on

both a and the actual situation in which the decision is taken. However, although the experimenter can quantify the consequences of his decisions for each possible situation, he is not sure what the situation actually is. So, he first must spend some time to collect information about the situation before committing himself to a specific decision. The only difficulty is that the information is not free. Thus, the experimenter must decide at which point he must *stop* collecting data and finally make a decision.

Such *optimal stopping* problems arise in many settings: Clinical trials Chernoff [1966], optimization Boender and Rinnooy Kan [1987], detecting changes in distributions Moustakides [1986], active learning Dimitrakakis and Savu-Krohn [2008]; Roy and McCallum [2001], as well as the the problem of deciding when to halt an optimization procedure Boender and Rinnooy Kan [1987]. A general introduction to such problems can be found in DeGroot [1970].

As before, we assume the existence of a *utility function* $U(a, \mu)$ defined for all $\mu \in \mathcal{M}$, where \mathcal{M} is the set of all possible universes of interest. The experimenter knows the utility function, but is not sure which universe the experiment is taking place in. This uncertainty about which $\mu \in \mathcal{M}$ is true is expressed via a subjective distribution $\xi_t(\mu) \triangleq \mathbf{P}(\mu|\xi_t)$, where ξ_t represents the belief at time t .

The expected utility of *immediately* taking an action at time t can then be written as $V_0(\xi_t) = \max_a \sum_{\mu} U(a, \mu) \xi_t(\mu)$, i.e. the experimenter takes the action which seems best on average. Now, consider that instead of making an immediate decision, he has the opportunity to take k more observations $D^k = (d_1, \dots, d_k)$ from a sample space S^k , at a cost $c > 0$ per observation⁴, thus allowing him to update his belief to

$$\xi_{t+k}(\mu|\xi_t) \triangleq \xi_t(\mu|D_k).$$

What the experimenter must do in order to choose between immediately making a decision a and continuing sampling, is to compare the utility of making a decision now with the cost of making k observations plus the utility of making a decision after k time-steps, when the extra data would enable a more informed choice.

The problem is in fact a dynamic programming problem. The utility of making an immediate decision is

$$V_0(\xi_t) = \max_a \int_{\mathcal{M}} U(a, \mu) \xi_t(\mu) d\mu \quad (3.1)$$

Similarly, we denote the utility of an immediate decision at any time $t + T$ by $V_0(\xi_{t+T})$. The utility of taking at most k samples before making a decision can be written recursively as:

$$V_{k+1}(\xi_t) = \max\{V_0(\xi_t), \mathbf{E}[V_k(\xi_{t+1})|\xi_t] - c\}, \quad (3.2)$$

where the expectation with respect to ξ_t is in fact taken over all possible observations under belief ξ_t :

$$\mathbf{E}[V_k(\xi_{t+1})|\xi_t] = \sum_{d_{t+1} \in S} V_k(\xi_{t+1}(\mu|d_{t+1})) \xi_t(d_{t+1}), \quad (3.3)$$

$$\xi_t(d_{t+1}) = \int_{\mathcal{M}} \mu(d_{t+1}) \xi_t(\mu) d\mu, \quad (3.4)$$

where $\xi_t(\mu|d_{t+1})$ indicates the specific next belief ξ_{t+1} arising from the previous belief ξ_t and the observations d_{t+1} .

This indicates that we can perform a backwards induction procedure, starting from all possible terminal belief states ξ_{t+T} to calculate the value of stopping immediately. We would only be able to insert the payoff function U directly when

⁴The case of non-constant cost is not significantly different.

we calculate V_0 . Taking $T \rightarrow \infty$ gives us the exact solution. In other words, one should stop and make an immediate decision if the following holds for all $k > 0$:

$$V_0(\xi_t) \geq V_k(\xi_t). \quad (3.5)$$

Note that if the payoff function is bounded we can stop the procedure at $T \propto c^{-1}$.

A number of bounds may be useful for stopping problems. For instance, the expected value of perfect information to an experimenter, can be used as a surrogate to the complete problem, and was examined by McCall [1965]. Conversely, lower bounds on the expected sample size and the regret were examined by Hoeffding [1960]. Stopping problems also appear in the context of bandit problems, or more generally, reinforcement learning. For instance the problem of finding a nearly optimal plan with high probability can be converted to a stopping problem Even-Dar et al. [2006].

3.2 Bandit problems

We now consider a generalization of the stopping problem. Imagine that our experimenter visits a casino and is faced with n different bandit machines. Playing a machine at time t results in a random payoff $r_t \in R \subset \mathbb{R}$. The average payoff of the i -th machine is μ_i . Thus, at time t the experimenter selects a machine with index $i \in \{1, \dots, n\}$ to receive a random payoff with expected value $\mathbf{E}(r_t | a_t = i) = \mu_i$, which is *fixed but unknown*. The experimenter's goal is to leave the casino with as much money as possible. This can be formalized as maximizing the expected sum of discounted future payoffs to time T :

$$\mathbf{E} \left(\sum_{t=1}^T \gamma^k r_t \right), \quad (3.6)$$

where $\gamma \in [0, 1]$ is a discount factor that reduces the importance of payoffs far in the future as it approaches 0. The horizon T may be finite or infinite, fixed, drawn from a known distribution, or simply unknown. It is also possible that the experimenter is allowed to stop at any time, thus adding stopping to the set of possible decisions and making the problem a straightforward generalization of the stopping problem. If the average payoffs are known, then the optimal solution is obviously to always take action $a^* = \arg \max_i \mu_i$, no matter what γ and T are and is thus completely uninteresting.

This problem was typically studied in a Bayesian setting Chernoff [1966]; Gittins [1989], where computable optimal solutions have been found for some special cases Gittins [1989]. However, recently there have been algorithms that achieve optimal regret rates in a distribution-free setting. In particular, the UCB1 algorithm by Auer et al [2002] selects the arm with highest empirical mean plus an upper confidence bound, with an error probability schedule tuned to achieve low regret. More specifically, let the empirical mean of the i -th arm at time t be:

$$\hat{\mathbf{E}}[r_t | a_t = i] \triangleq \frac{1}{n_i^t} \sum_{k: a_k = i}^t r_k, \quad n_i^t \triangleq \sum_{k=1}^t \mathbb{I}\{a_k = i\}, \quad (3.7)$$

where n_i^t is the number of times arm i has been played until time t . After playing each arm once, the algorithm always selects the arm maximizing:

$$\hat{\mathbf{E}}[r_t | a_t = i] + \sqrt{\frac{2 \log t}{n_i^t}}. \quad (3.8)$$

This guarantees a regret that only scales with rate $O(\log T)$.

The setting has been generalized to continuous time Chernoff [1966], non-stationary or adversarial bandits Auer [2002], continuous spaces Agrawal [1995]; Auer et al. [2007] and to trees Coquelin and Munos [2007]; Kocsis and Szepesvári [2006], while a collection of related results can be found in [Cesa-Bianchi and Lugosi, 2006]. Finally, the bandit payoffs can also depend on a context variable. If this variable cannot be affected by the experimenter, then it is sufficient to learn the mean payoffs for all contexts in order to be optimal. However, the problem becomes much more interesting when the experimenter's actions also influence the context of the game. This directly leads us to the concept of problems with state.

3.3 Reinforcement learning and control

We now generalize further to problems where the payoffs depend not only on the individual actions that we perform, but also on a context, or state. This is a common situation in games. A good example is blackjack, where drawing or stopping (the two possible actions in a game) have expected payoffs that depend on your current hand (the current state, assuming the croupier's hand is random and independent of your own hand).

Both reinforcement learning and control problems are formally identical. Nevertheless, historically, classical control (c.f. Stengel [1994]) addressed the case where the objective is a known functional of the state s and action a . Reinforcement learning, on the other hand, started from the assumption that the objective function itself is unknown (though its functional *form* is known) and must be estimated. Both discrete-time control and reinforcement learning problems can be formalized in the framework of Markov decision processes.

4 Markov decision processes

Definition 4.1 (Markov decision process) A Markov decision process (MDP) is defined as the tuple $\mu = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R})$ comprised of a set of states \mathcal{S} , a set of actions \mathcal{A} , a transition distribution \mathcal{T} conditioning the next state on the current state and action,

$$\mathcal{T}(s' | s, a) \triangleq \mu(s_{t+1}=s' | s_t=s, a_t=a) \quad (4.1)$$

satisfying the Markov property $\mu(s_{t+1} | s_t, a_t) = \mu(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, \dots)$, and a reward distribution \mathcal{R} conditioned on states and actions:

$$\mathcal{R}(r | s, a) \triangleq \mu(r_{t+1}=r | s_t=s, a_t=a), \quad (4.2)$$

with $a \in \mathcal{A}$, $s, s' \in \mathcal{S}$, $r \in \mathbb{R}$. Finally,

$$\mu(r_{t+1}, s_{t+1} | s_t, a_t) = \mu(r_{t+1} | s_t, a_t) \mu(s_{t+1} | s_t, a_t). \quad (4.3)$$

We shall denote the set of all MDPs as \mathcal{M} . We are interested in sequential decision problems where, at each time step t , the agent seeks to maximize the expected utility

$$\sum_{k=1}^{T-t} \gamma^k \mathbf{E}[r_{t+k} | \cdot],$$

where r is a stochastic reward and u_t is simply the discounted sum of future rewards. We shall assume that the sequence of rewards arises from a Markov decision process.

In order to describe how we select actions, we define a policy π as a distribution on actions conditioned on the current state $\pi(a_t|s_t)$. If we are interested in the rewards to some horizon T , then we can define a T -horizon value function for an MDP $\mu \in \mathcal{M}$ at time t as:

$$V_{\mu,t,T}^{\pi}(s) \triangleq \sum_{k=1}^{T-t} \gamma^k \mathbf{E}[r_{t+k} | s_t = s, \pi, \mu], \quad (4.4)$$

the expected sum of future rewards given that we are at state s at time t and selecting actions according to policy π in the MDP μ . Superscripts and subscripts of V may be dropped when they are clear from context.

The value function of any policy can be calculated recursively by starting from the terminal states.

$$V_{\mu,t,T}^{\pi}(s) = \mathbf{E}[r_{t+1} | s_t = s, \pi, \mu] + \gamma \sum_{s'} \mu(s_{t+1} = s' | s_t = s, a_t = a) V_{\mu,t+1,T}^{\pi}(s'). \quad (4.5)$$

The optimal value function, i.e. the value function of the optimal policy π^* , can be calculated by a maximizing over actions at each stage:

$$V_{\mu,t,T}^{\pi^*}(s) = \max_{a \in \mathcal{A}} \mathbf{E}[r_{t+1} | s_t = s, a_t = a, \mu] + \gamma \sum_{s'} \mu(s_{t+1} = s' | s_t = s, a_t = a) V_{\mu,t+1,T}^{\pi^*}(s'). \quad (4.6)$$

The recursive form of this equation is frequently referred to as the Bellman recursion and allows us to compute the value of a predecessor state from that of a following state. The resulting algorithm is called *backwards induction* or *value iteration*. When the number of states is finite, the same recursion allows us to compute the value of states when the horizon is infinite, as then, $\lim_{T \rightarrow \infty} V_{\mu,t,T}^{\pi} = V_{\mu}^{\pi}$ for all finite t . Finally note that frequently we shall denote $V_{\mu,t,T}^{\pi^*}$ simply by V^* when the environmental variables are clear from context.

5 Belief-augmented Markov decision processes

When the MDP is unknown, we need to explicitly take into account our uncertainty. In control theory, this is referred to as the problem of dual control [Stengel, 1994, Sec. 5.2]. This involves selecting actions (control inputs) such as to improve parameter (and state) estimation in order to hopefully reduce future costs. This behavior is called probing. At the same time, the control strategy must not neglect the minimization of the cost at the current time.

This type of dilemma in its simplest forms occurs in the already discussed bandit problems, where we must strike an optimal balance between exploring alternative bandits and exploiting the apparently best bandit. Any optimal solution must take into account the uncertainty that we have about the environment.

A natural idea is to use a Bayesian framework (c.f. Duff [2002]) to represent our beliefs. As summarized in section 2.3, this involves maintaining a belief $\xi_t \in \Xi$, about which MDP $\mu \in \mathcal{M}$ corresponds to reality. In a Bayesian setting, $\xi_t(\mu)$ is a subjective probability density over MDPs.

5.1 Bayesian inference with a single MDP model class

We shall cover the case where it is known that the true MDP μ^* is in some set of MDPs \mathcal{M} . For example, it may be known that the MDP has at most K discrete states and that the rewards at each state are Bernoulli, but we know neither the actual transition probabilities nor the reward distributions. Nevertheless, we can use

closed form Bayesian techniques to model our belief about which model is correct: For discrete state spaces, transitions can be expressed as multinomial distributions, to which the Dirichlet density is a conjugate prior. Similarly, unknown Bernoulli distributions can be modelled via a Beta prior. If we assume that the densities are not dependent, then the prior over all densities is a product of Dirichlet priors, and similarly for rewards. Then we only need a number of parameters of order $O(|\mathcal{S}|^2|\mathcal{A}|)$. The remainder of this section discusses this in more detail.

Let \mathcal{M} be the set of MDPs with unknown transition probabilities and state space \mathcal{S} of size K . We denote our belief at time $t + 1$ about which MDP is true as simply our belief density at time t conditioned on the latest observations:

$$\xi_{t+1}(\mu) \triangleq \xi_t(\mu | r_{t+1}, s_{t+1}, s_t, a_t) \quad (5.1a)$$

$$= \frac{\mu(r_{t+1}, s_{t+1} | s_t, a_t) \xi_t(\mu)}{\int_{\mathcal{M}} \mu'(r_{t+1}, s_{t+1} | s_t, a_t) \xi_t(\mu') d\mu'}. \quad (5.1b)$$

We consider the case where \mathcal{M} is an infinite set of MDPs, where each MDP $\mu \in \mathcal{M}$ corresponds to a particular joint probability distribution over the state-action pairs.

We shall begin by defining a belief for the transition of each state action pair s, a separately. Firstly, we denote by $\tau_{s,a} \in [0, 1]^K$ the parameters of the multinomial distribution over the K possible next states, from a specific starting state s and action a . Our belief will be a Dirichlet distribution – a function of $x \in \mathbb{R}^K$ with $\|x\|_1 = 1$ and $x \in [0, 1]^K$, with parameters $\psi^{s,a} \in \mathbb{N}^K$. If we denote the parameters of our belief ξ_t at time t by $\psi^{s,a}(t)$, then the Dirichlet density over possible multinomial distributions can be written as:

$$\xi_t(\tau_{s,a} = x) = \frac{\Gamma(\psi^{s,a}(t))}{\prod_{i \in \mathcal{S}} \Gamma(\psi_i^{s,a}(t))} \prod_{i \in \mathcal{S}} x_i^{\psi_i^{s,a}(t)}, \quad (5.2)$$

where $\psi_i^{s,a}$ denotes the i -th component of $\psi^{s,a}$. The set of parameters ψ can be written in matrix form as $\Psi(t)$ to denote the $|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|$ matrix of state-action-state transition counts at time t . The *initial* parameters $\Psi(0)$ form a matrix that defines the parameters of our set *prior* Dirichlet distributions.

Thus, for any belief ξ_t , the Dirichlet parameters are $\{\psi_i^{j,a}(t) : i, j \in \mathcal{S}, a \in \mathcal{A}\}$. These values are initialised to $\Psi(0)$ and are updated via simple counting:

$$\psi_i^{j,a}(t+1) = \psi_i^{j,a}(t) + \mathbb{I}\{s_{t+1} = i \wedge s_t = j \wedge a_t = a\}, \quad (5.3)$$

meaning that every time we observe a specific transition s_t, a_t, s_{t+1} , we increment the corresponding Dirichlet parameter by one.

We now need to move from the distribution of a single state-action pair to the set of transition distributions for the whole MDP. In order to do this easily, we shall make the following simplifying assumption:

Assumption 5.1 For any $s, s' \in \mathcal{S}$ and $a, a' \in \mathcal{A}$,

$$\xi(\tau_{s,a}, \tau_{s',a'}) = \xi(\tau_{s,a}) \xi(\tau_{s',a'}). \quad (5.4)$$

This assumption significantly simplifies the model but does not let us take into advantage of the case where there may be some dependencies in the transition probabilities. Now we shall denote the matrix of state-action-state transition *probabilities* for a specific MDP μ as \mathcal{T}^μ . Analogously to $\tau_{s,a}$, we denote, for the

specific MDP μ , the next state distribution multinomial parameter vector from pair (s, a) to be $\tau_{s,a}^\mu$, with $\tau_{s,a}^\mu(i) \triangleq \mu(s_{t+1} = i \mid s_t = s, a_t = a)$. Then we obtain:

$$\xi_t(\mu) = \xi_t(\mathcal{T}^\mu) = \xi_t(\tau_{s,a} = \tau_{s,a}^\mu \forall s \in \mathcal{S}, a \in \mathcal{A}) \quad (5.5a)$$

$$= \prod_{s \in \mathcal{S}} \prod_{a \in \mathcal{A}} \xi_t(\tau_{s,a} = \tau_{s,a}^\mu), \quad (5.5b)$$

$$= \prod_{s \in \mathcal{S}} \prod_{a \in \mathcal{A}} \frac{\Gamma(\Psi^{s,a}(t))}{\prod_{i \in \mathcal{S}} \Gamma(\Psi_i^{s,a}(t))} \prod_{i \in \mathcal{S}} \left(\tau_{s,a,i}^\mu \right)^{\Psi_i^{s,a}(t)}, \quad (5.5c)$$

where we used Assumption 5.1. This means that the transition counts Ψ are a sufficient statistic for expressing the density over \mathcal{M} .

We can additionally model $\mu(r_{t+1} \mid s_t, a_t)$ with a suitable belief (for example a Beta, Normal-Gamma or Pareto prior) and assume independence. This in no way complicates the exposition for MDPs.

5.2 Constructing and solving BAMDPs

In order to optimally select actions in this framework, we need to consider how our actions will affect our future beliefs. The corresponding Bayesian procedure is not substantially different from the optimal stopping procedure outlined in 3.1. The approach outlined in this section was suggested originally in Bellman and Kalaba [1959] under the name of Adaptive Control Processes and was investigated more fully in Duff and Barto [1997]; Duff [2002]. It involves the creation of an *augmented* MDP, with a state comprised of the original MDP's state s_t and our belief state ξ_t . We can then solve the problem *in principle* via standard dynamic programming algorithms such as backwards induction, similarly to section 3.1. We shall call such models BAMDPs (Belief-Augmented MDPs).

In BAMDPs, we are at some combined belief and environment state $\omega_t = (\xi_t, s_t)$ at each point in time t , which we call the *hyper-state*. For every possible action a_t , we may observe any $s_{t+1} \in \mathcal{S}$ and any possible reward $r_{t+1} \in R \subset \mathbb{R}$, which would lead to a unique new belief ξ_{t+1} and thus a unique new hyper-state $\omega_{t+1} = (\xi_{t+1}, s_{t+1})$.

More formally, we may give the following definition:

Definition 5.1 (Belief-Augmented MDP) A Belief-Augmented MDP v (BAMDP) is an MDP $v = (\Omega, \mathcal{A}, \mathcal{T}', \mathcal{R}')$ where $\Omega = \mathcal{S} \times \Xi$, where Ξ is an appropriate set of probability measures on \mathcal{M} , and $\mathcal{T}', \mathcal{R}'$ are the transition and reward distributions conditioned jointly on the MDP state s_t , the belief state ξ_t , and the action a_t . Here the density $p(\xi_{t+1} \mid \xi_t, r_{t+1}, s_{t+1}, s_t, a_t)$ is singular, since ξ_{t+1} is a deterministic function of $\xi_t, r_{t+1}, s_{t+1}, s_t, a_t$. Thus, we can define the transition

$$v(\omega_{t+1} \mid a_t, \omega_t), \quad (5.6)$$

where $\omega_t \triangleq (s_t, \xi_t)$.

It should be obvious that s_t, ξ_t jointly form a Markov state in this setting, called the *hyper-state*. In general, we shall denote the components of a future hyper-state ω_t^i as (s_t^i, ξ_t^i) . However, in occasion we will abuse notation by referring to the components of some hyper-state ω as s_ω, ξ_ω . We shall use \mathcal{M}_B to denote the set of BAMDPs.

As in the MDP case, finite horizon problems only require looking at all future hyper-states until the horizon T , where we omit the subscript v for the value function:

$$V_{t,T}^*(\omega_t) = \max_{a_t} \mathbf{E}[r_{t+1} \mid \omega_t, a_t, v] + \gamma \int_{\Omega} V_{t+1,T}^*(\omega_{t+1}) v(\omega_{t+1} \mid \omega_t, a_t) d\omega_{t+1}. \quad (5.7)$$

It is easy to see that the size of the set of hyper-states in general grows exponentially with the horizon. Thus, we can not perform value iteration with bounded memory, as we did for discrete MDPs. One possibility is to continue expanding the belief tree until we are certain of the optimality of an action. As has previously been observed Dearden et al. [1998]; Dimitrakakis [2006], this is possible since we can always obtain upper and lower bounds on the utility of any policy from the current hyper-state. In addition, we can apply such bounds on future hyper-states in order to efficiently expand the tree.

5.3 Belief tree expansion

Let the current belief be ξ_t and suppose we observe $x_t^i \triangleq (s_{t+1}^i, r_{t+1}^i, a_t^i)$. This observation defines a unique subsequent belief ξ_{t+1}^i . Together with the MDP state s , this creates a hyper-state transition from ω_t to ω_{t+1}^i .

By recursively obtaining observations for future beliefs, we can obtain an unbalanced tree with nodes $\{\omega_{t+k}^i : k = 1, \dots, T; i = 1, \dots\}$. However, we cannot hope to be able to fully expand the tree. This is especially true in the case where observations (i.e. states, rewards, or actions) are continuous, where we cannot perform even a full single-step expansion. Even in the discrete case the problem is intractable for infinite horizons – and far too complex computationally for the finite horizon case. However, had there been efficient tree expansion methods, this problem would be largely alleviated.

All tree search methods require the expansion of leaf nodes. However, in general, a leaf node may have an infinite number of children. We thus need some strategies to limit the number of children. More formally, let us assume that we wish to expand in node $\omega_t^i = (\xi_t^i, s_t^i)$, with ξ_t^i defining a density over \mathcal{M} . For discrete state/action/reward spaces, we can simply enumerate all the possible outcomes $\{\omega_{t+1}^j\}_{j=1}^{|\mathcal{S} \times \mathcal{A} \times R|}$, where R is the set of possible reward outcomes. Note that if the reward is deterministic, there is only one possible outcome per state-action pair. The same holds if \mathcal{S} is deterministic, in both cases making an enumeration possible. While in general this may not be the case, since rewards, states, or actions can be continuous, in this chapter we shall only examine the discrete case.

5.4 Bounds on the optimal value function

Let Ω_T be the set of leaf nodes of the partially expanded belief tree and v the BAMDP process. If the values of the leaf nodes were known, then we could easily perform the backwards induction procedure, shown in Algorithm 1 for BAMDPs: If one thinks of the BAMDP as a very large MDP, one can see that the algorithm (also called value iteration) is identical to equation (4.6), with the subtle difference that the reward only depends on the next hyper-state.

The main problem is obtaining a good estimate for V_T^* , i.e. the value of leaf nodes. Let $\pi^*(\mu)$ denote the policy such that, for any π ,

$$V_\mu^{\pi^*(\mu)}(s) \geq V_\mu^\pi(s), \quad \forall s \in \mathcal{S}.$$

Furthermore, let the mean MDP arising from the belief ξ at hyper-state $\omega = (s, \xi)$ be $\bar{\mu}_\xi \triangleq \mathbf{E}[\mu|\xi]$.

Proposition 5.1 *The optimal value function V^* of the BAMDP v at any hyper-state $\omega = (s, \xi)$ is bounded by the following inequalities*

$$\int V_\mu^{\pi^*(\mu)}(s) \xi(\mu) d\mu \geq V^*(\omega) \geq \int V_\mu^{\pi^*(\bar{\mu}_\xi)}(s) \xi(\mu) d\mu. \quad (5.8)$$

Algorithm 1 Backwards induction action selection

1: **input** Process ν , time t , leaf nodes Ω_T , leaf node values V_T^* .
2: **for** $n = T - 1, T - 2, \dots, t$ **do**
3: **for** $\omega \in \Omega_n$ **do**
4:

$$a_n^*(\omega) = \arg \max_a \sum_{\omega' \in \Omega_{n+1}} \nu(\omega' | \omega, a) [\mathbf{E}(r | \omega', \omega, \nu) + V_{n+1}^*(\omega')]$$
$$V_n^*(\omega) = \sum_{\omega' \in \Omega_{n+1}} \nu(\omega' | \omega, a^*) [\mathbf{E}(r | \omega', \omega, \nu) + V_{n+1}^*(\omega')]$$

5: **end for**
6: **end for**
7: **return** a_t^*

The proof is given in the appendix. In POMDPs, a trivial lower bound can be obtained by calculating the value of the blind policy Hauskrecht [2000]; Smith and Simmons [2005], which always takes the a *fixed* action, i.e. $a_t = a$ for all t . Our lower bound is in fact the BAMDP analogue of the value of the blind policy in POMDPs if we consider BAMDP policies are POMDP actions. The analogy is due to the fact that for any policy π , which selects actions by considering only the MDP state, i.e. such that $\pi(a_t | s_t, \xi_t) = \pi(a_t | s_t)$, it holds trivially that $V^\pi(\omega) \leq V^*(\omega)$, in the same way that it holds trivially if we consider only the set of policies which always take the same action. In fact, of course $V^*(\omega) \geq V^\pi(\omega)$ for any π , by definition. In our case, we have made this lower bound tighter by considering $\pi^*(\bar{\mu}_\xi)$, the policy that is greedy with respect to the current mean estimate.

The upper bound itself is analogous to the POMDP value function bound given in Theorem 9 of Hauskrecht [2000]. The crucial difference is that, in our case, both bounds can only be approximated via Monte Carlo sampling with some probability, unless \mathcal{M} is finite.

5.4.1 Leaf node lower bound

A lower bound can be obtained by calculating the expected value of any policy. In order to have a tighter bound, we can perform value iteration in the mean MDP. Let us use $\bar{\mu}_\xi$ to denote the *mean MDP* for belief ξ , with transition probabilities $\mathcal{T}_{\bar{\mu}_\xi}$ and mean rewards $R_{\bar{\mu}_\xi}$:

$$\mathcal{T}_{\bar{\mu}_\xi} \triangleq \bar{\mu}_\xi(s_{t+1} | s_t, a_t) = \mathbf{E}(\mathcal{T}^\mu | \xi_t) \quad (5.9)$$

$$R_{\bar{\mu}_\xi} \triangleq \bar{\mu}_\xi(s_{t+1} | s_t, a_t) = \mathbf{E}(\mathcal{R}^\mu | \xi_t). \quad (5.10)$$

Similarly, let $V_{\bar{\mu}_\xi}^\pi$ be the column vector of the value function of the mean MDP, to obtain:

$$\begin{aligned} V_{\bar{\mu}_\xi}^\pi &= R_{\bar{\mu}_\xi} + \gamma \int \mathcal{T}_{\bar{\mu}_\xi}^\pi V_{\bar{\mu}_\xi}^\pi \xi(\mu) d\mu \\ &= R_{\bar{\mu}_\xi} + \gamma \left(\int \mathcal{T}_{\bar{\mu}_\xi}^\pi \xi(\mu) d\mu \right) V_{\bar{\mu}_\xi}^\pi \\ &= R_{\bar{\mu}_\xi} + \gamma \mathcal{T}_{\bar{\mu}_\xi}^\pi V_{\bar{\mu}_\xi}^\pi. \end{aligned}$$

This is now a standard Bellman recursion, which we can use to obtain the policy $\pi_{\bar{\mu}_\xi}^*$ which is optimal with respect to the mean MDP. Unfortunately the value function

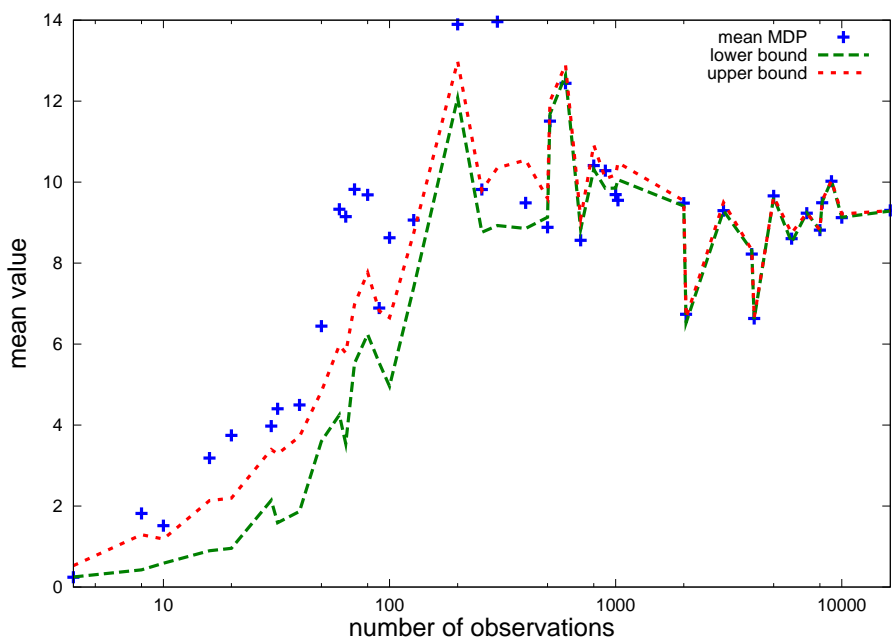


Figure 1: Illustration of upper and lower bounds on the value function, averaged over all states, as more observations are acquired. It can be seen that the , mean MDP value (crosses) is far from the bounds initially. The bounds were calculated by taking the empirical mean of 1000 MDP samples from the current belief.

of the mean MDP does not generally equal the expected value of the BAMDP:

$$V_{\bar{\mu}_\xi}^\pi \neq \mathbf{E}[V^\pi | \xi].$$

Nevertheless, the stationary policy that is optimal with respect to the mean MDP can be used to evaluate the right hand side for $\pi = \pi_{\bar{\mu}}^*$, which will hopefully result in a relatively tight lower bound. This is illustrated in Figure 1, which displays upper and lower bounds on the BAMDP value function as observations are acquired in a simple maze task.

If the beliefs ξ can be expressed in closed form, it is easy to calculate the mean transition distribution and the mean reward from ξ . For discrete state spaces, transitions can be expressed as multinomial distributions, to which the Dirichlet density is a conjugate prior. In that case, for Dirichlet parameters $\{\psi_i^{j,a}(\xi) : i, j \in \mathcal{S}, a \in \mathcal{A}\}$, we have

$$\bar{\mu}_\xi(s' | s, a) = \frac{\psi_{s'}^{s,a}(\xi)}{\sum_{i \in \mathcal{S}} \psi_i^{s,a}(\xi)} \quad (5.11)$$

Similarly, for Bernoulli rewards, the corresponding mean model arising from the beta prior with parameters $\{\alpha^{s,a}(\xi), \beta^{s,a}(\xi) : s \in \mathcal{S}, a \in \mathcal{A}\}$ is $\mathbf{E}[r | s, a, \bar{\mu}_\xi] = \alpha^{s,a}(\xi) / (\alpha^{s,a}(\xi) + \beta^{s,a}(\xi))$. Then the value function of the mean model can be found with standard value iteration.

5.4.2 Bounds with high probability

In general, neither the upper nor the lower bounds cannot be expressed in closed form. However, the integral can be approximated via Monte Carlo sampling.

Let us be in some hyper-state $\omega = (s, \xi)$. We can obtain c MDP samples from the belief at ω : $\mu_1, \dots, \mu_c \sim \xi(\mu)$. In order to estimate the upper bound, for each μ_k we can derive the optimal policy $\pi^*(\mu_k)$ and estimate its value function $\tilde{v}_k^* \triangleq V_{\mu_k}^{\pi^*(\mu_k)} \equiv V_{\mu_k}^*$. We may then average these samples to obtain

$$\hat{v}_c^*(\omega) \triangleq \frac{1}{c} \sum_{k=1}^c \tilde{v}_k^*(s), \quad (5.12)$$

where s is the state at hyper-state ω . Let $\bar{v}^*(\omega) = \int_{\mathcal{M}} \xi(\mu) V_\mu^*(s) d\mu$. It holds that $\lim_{c \rightarrow \infty} [\hat{v}_c^*] = \bar{v}^*(\omega)$ and that $\mathbf{E}[\hat{v}_c^*] = \bar{v}^*(\omega)$. Due to the latter, we can apply a Hoeffding inequality

$$\mathbf{P}(|\hat{v}_c^*(\omega) - \bar{v}^*(\omega)| > \varepsilon) < 2 \exp\left(-\frac{2c\varepsilon^2}{(V_{\max} - V_{\min})^2}\right), \quad (5.13)$$

thus bounding the error within which we estimate the upper bound. For $r_t \in [0, 1]$ and discount factor γ , note that $V_{\max} - V_{\min} \leq 1/(1 - \gamma)$.

The procedure for the lower bound is identical, but we only need to estimate the value $\tilde{v}_k \triangleq V_{\mu_k}^{\pi^*(\bar{\mu}_\xi)}$ of the mean-MDP-optimal policy $\pi^*(\bar{\mu}_\xi)$ for each one of the sampled MDPs μ_k . Thus we obtain a pair of high probability bounds for any BAMDP node.

5.5 Discussion and related work

Employing the above bounds together with efficient search methods [Coquelin and Munos, 2007; Dimitrakakis, 2008, 2009; Hren and Munos, 2008] is a promising direction. Even when the search is efficient, however, the complexity remains prohibitive for large problems due to the high branching factor and the dependency on the horizon.

Poupart et al [Poupart et al. [2006]] have proposed an analytic solution to Bayesian reinforcement learning. They focus on how to efficiently approximate (5.7). More specifically, they use the fact that the optimal value function is the upper envelope of a set of linear segments

$$V_{t,T}^*(\omega_t) = \max_{\alpha \in \Gamma} \alpha(\omega_t),$$

with

$$\alpha(\omega_t) = \int_{\mathcal{M}} \xi_i(\mu) \alpha(s_t, \mu) d\mu.$$

They then show that the $k+1$ -horizon α -function can be computed from the k -horizon α -function via the following backwards induction:

$$\alpha(\omega_t) = \sum_i \mu(s_{t+1}=i | s_t, a^*(\omega_t)) [\mathbf{E}(R | s_{t+1}=i, s_t, a^*(\omega_t)) + \gamma \alpha(\omega_{t+1}^{i, a^*(\omega)})], \quad (5.14)$$

$$V^{k+1}(\omega) = \max_{\alpha \in \Gamma^{k+1}} \alpha(\omega), \quad (5.15)$$

where ω_{t+1}^i denotes the augmented state resulting from starting at ω_t , taking action $a^*(\omega_t)$ and transiting to the MDP state i . However, the complexity of this process is still exponential with the planning horizon. For this reason, the authors use a projection of the α -functions. This is performed on a set of belief points selected via sampling a default trajectory. The idea is that one may generalize from the set of sample beliefs to other, similar, belief states. In fact, the approximate value function they derive is a lower bound on the BAMDP value function. It is an open question whether or when the bounds presented here are tighter.

At this point, however, it is unclear under what conditions efficient online methods would outperform approximate analytic methods. Perhaps a combination of the two methods would be of interest: i.e. using the online search and bounds in order to see when the point-based approximation has become too inaccurate. The online search could also be used to sample new belief points.

It should be noted that online methods break down when $\gamma \rightarrow 1$ because the belief tree can not in general be expanded to the required depth. In fact, the only currently known methods which are nearly optimal in the undiscounted case are based on distribution-free bounds Auer et al. [2008]. Similarly to the bandit case, it is possible to perform nearly optimally in an unknown MDP by considering upper confidence bounds on the value function. Auer et al [Auer et al., 2008] in fact give an algorithm which achieves regret $O(D|\mathcal{S}|\sqrt{|\mathcal{A}|T})$ after T steps for any unknown MDP with diameter⁵ D . In exactly the same way as the Bayesian approach, the algorithm maintains counts Ψ over observe state-action-state transitions. In order to obtain an optimistic value function, the authors consider an augmented MDP which is constructed from a set of plausible MDPs M , where M is such that $\mathbf{P}(\Psi | \mu \notin M) < \delta$. Then, instead of augmenting the state space, they augment the action space by allowing the simultaneous choice of actions $a \in \mathcal{A}$ and MDPs $\mu \in M$. The policy is then chosen by performing average value iteration [c.f. Puterman, 2005] in the augmented MDP.

There has not been much work yet for Bayesian methods in the undiscounted case, with the exception of Bernoulli bandit problems Kelly [1981]. It may well be that in fact naive look-ahead methods are impractical. In order to achieve $O(D|\mathcal{S}|\sqrt{|\mathcal{A}|T})$ regret with online methods Dimitrakakis [2009] requires an instantaneous regret ε_t such that $\sum_t^T \varepsilon_t < \sqrt{T}$, $\Rightarrow \varepsilon_t < 1/\sqrt{t}$. If we naively bound our instantaneous regret by discounting, then that would require our horizon to grow

⁵Intuitively, the maximin expected time needed to reach any state from any other state, where the max is taken over state pairs and the min over policies. See [Puterman, 2005] for details.

with rate $O(\log_{1/\gamma} \sqrt{t}/(1-\gamma))$, something which seems impractical with current online search techniques.

6 Partial observability

A useful extension of the MDP model can be obtained by not allowing the agent to directly observe the state of the environment, but an observation variable o_t that is conditioned on the state. This more realistic assumption is formally defined as follows:

Definition 6.1 (Partially observable Markov decision process) A *partially observable Markov decision process* μ (POMDP) is defined as the tuple $\mu = (\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R})$ comprised of a set of observations \mathcal{O} , a set of states \mathcal{S} , a set of actions \mathcal{A} , a transition-observation distribution \mathcal{T} conditioned the current state and action $\mu(s_{t+1} = s', o_{t+1} = o | s_t = s, a_t = a)$ and a reward distribution \mathcal{R} , conditioned on the state and action $\mu(r_{t+1} = r | s_t = s, a_t = a)$, with $a \in \mathcal{A}$, $s, s' \in \mathcal{S}$, $o \in \mathcal{O}$, $r \in \mathbb{R}$.

We shall denote the set of POMDPs as \mathcal{M}_P . For POMDPs, it is often assumed that one of the two following factorizations holds:

$$\mu(s_{t+1}, o_{t+1} | s_t, a_t) = \mu(s_{t+1} | s_t, a_t) \mu(o_{t+1} | s_{t+1}) \quad (6.1)$$

$$\mu(s_{t+1}, o_{t+1} | s_t, a_t) = \mu(s_{t+1} | s_t, a_t) \mu(o_{t+1} | s_t, a_t). \quad (6.2)$$

The assumption that the observations are only dependent on a single state or a single state-action pair is a natural decomposition for a lot of practical problems.

POMDPs are *formally identical* to BAMDPs. More specifically, BAMDPs correspond to a special case of a POMDP in which the state is split into two parts: One fully observable dynamic part and one unobservable, continuous, but stationary part, which models the unknown MDP. Typically, however, in POMDP applications the unobserved part of a state is dynamic and discrete.

The problem of acting optimally in POMDPs has two aspects. The first is state estimation, and the second is acting optimally given the estimated state. As far as the first part is concerned, given an initial state probability distribution, updating the belief for a discrete state space amounts to simply maintaining a multinomial distribution over the states. However, the initial state distribution might not be known. In that case, we may assume an initial prior density over the multinomial state distribution. It is easy to see that this is simply a special case of an unknown state transition distribution, where we insert a special initial state which is only visited once. We shall, however, be concerned with the more general case of full exploration in POMDPs, where all state transition distributions are unknown.

6.1 Belief POMDPs

It is possible to create an augmented MDP for POMDP models, by endowing them with an additional belief state, in the same manner as MDPs. However now the belief state will be a joint probability distribution over \mathcal{M}_P and \mathcal{S} . Nevertheless, each (a_t, o_{t+1}) pair that is observed leads to a unique subsequent belief state. More formally, a belief-augmented POMDP is defined as follows:

Definition 6.2 (Belief POMDP) A Belief POMDP v (BAPOMPD) is an MDP $v = (\Omega, \mathcal{A}, \mathcal{O}, \mathcal{T}', \mathcal{R}')$ where $\Omega = \mathcal{G} \times \mathcal{B}$, where \mathcal{G} is the set of probability measures on \mathcal{S} , \mathcal{B} is the set of probability measures on \mathcal{M}_P , \mathcal{T}' \mathcal{R}' are the belief

state transition and reward distributions conditioned on the belief state ξ_t and the action a_t such that the following factorizations are satisfied for all $\mu \in \mathcal{M}_P$, $\xi_t \in \mathcal{B}$

$$p(s_{t+1}|s_t, a_t, s_{t-1}, \dots, \mu) = \mu(s_{t+1}|s_t, a_t) \quad (6.3)$$

$$p(o_t|s_t, a_t, o_{t-1}, \dots, \mu) = \mu(o_t|s_t, a_t) \quad (6.4)$$

$$p(\xi_{t+1}|o_{t+1}, a_t, \xi_t) = \int_{\mathcal{M}_P} p(\xi_{t+1}|\mu, o_{t+1}, a_t, \xi_t) \xi_{t+1}(\mu|o_{t+1}, a_t, \xi_t) d\mu \quad (6.5)$$

We shall denote the set of BAPOMDPs with \mathcal{M}_{BP} . Again, (6.5) simply assures that the transitions in the belief-POMDP are well-defined. The Markov state $\xi_t(\mu, s_t)$ now jointly specifies a distribution over POMDPs and states.⁶ As in the MDP case, in order to be able to evaluate policies and select actions optimally, we need to first construct the BAPOMDP. This requires calculating the transitions from the current belief state to subsequent ones according to our possible future observations, as well as the probability of those observations. The next section goes into this in more detail.

6.2 The belief state

In order to simplify the exposition, in the following we shall assume firstly that each POMDP has the same number of states. Then $\xi(s_t = s|\mu)$ describes the probability that we are in state s at time t given some belief ξ and assuming we are in the POMDP μ . Similarly, $\xi(s_t = s, \mu)$ is the joint probability given our belief. This joint distribution can be used as a state in an expanded MDP, which can be solved via backward induction, as will be seen later. In order to do this, we must start with an initial belief ξ_0 and calculate all possible subsequent beliefs. The belief at time $t+1$ depends only on the belief time t and the current set of observations r_{t+1}, o_{t+1}, a_t . Thus, the transition probability from ξ_t to ξ_{t+1} is just the probability of the observations according to our current belief, $\xi_t(r_{t+1}, o_{t+1}|a_t)$. This can be calculated by first noting that given the model and the state, the probability of the observations no longer depends on the belief, i.e.

$$\xi_t(r_{t+1}, o_{t+1}, |s_t, a_t, \mu) = \mu(r_{t+1}, o_{t+1} | a_t, s_t) = \mu(r_{t+1} | a_t, s_t) \mu(o_{t+1} | a_t, s_t). \quad (6.6)$$

The probability of any particular observation can be obtained by integrating over all the possible models and states

$$\xi_t(r_{t+1}, o_{t+1} | a_t) = \int_{\mathcal{M}_P} \int_{\mathcal{S}} \mu(r_{t+1}, o_{t+1} | a_t, s_t) \xi_t(\mu, s_t) d\mu ds_t. \quad (6.7)$$

Given that a particular observation is made from a specific belief state, we now need to calculate what belief state it would lead to. For this we need to compute the posterior belief over POMDPs and states. The belief over POMDPs is given by

$$\xi_{t+1}(\mu) \triangleq \xi_t(\mu | r_{t+1}, o_{t+1}, a_t) \quad (6.8)$$

$$= \frac{\xi_t(r_{t+1}, o_{t+1}, a_t | \mu) \xi_t(\mu)}{\xi_t(r_{t+1}, o_{t+1}, a_t)} \quad (6.9)$$

$$= \frac{\xi_t(\mu)}{Z} \iint_{\mathcal{S}} \mu(r_{t+1}, o_{t+1}, a_t | s_{t+1}, s_t) \xi_t(s_{t+1}, s_t | \mu) ds_{t+1} ds_t \quad (6.10)$$

⁶The formalism is very similar to that described in Ross et al. [2008a], with the exception that we do not include the actual POMDP state in the model.

where $Z = \xi_t(r_{t+1}, o_{t+1}, a_t)$ is a normalizing constant. Note that $\xi_t(s_{t+1}, s_t | \mu) = \mu(s_{t+1} | s_t) \xi_t(s_t | \mu)$, where $\xi_t(s_t | \mu)$ is our belief about the state in the POMDP μ . This can be updated using the following two steps. Firstly, the filtering step

$$\xi_{t+1}(s_t | \mu) \triangleq \xi_t(s_t | r_{t+1}, o_{t+1}, a_t, \mu) \quad (6.11)$$

$$= \frac{\mu(r_{t+1}, o_{t+1} | s_t, a_t) \xi_t(s_t | \mu)}{\xi_t(r_{t+1}, o_{t+1} | a_t, \mu)}, \quad (6.12)$$

where we adjust our belief about the previous state of the MDP based on. Then we must perform a prediction step

$$\xi_{t+1}(s_{t+1} | \mu) = \int_{\mathcal{S}} \mu(s_{t+1} | s) \xi_{t+1}(s_t = s | \mu) ds, \quad (6.13)$$

where we calculate the probability over the current states given our new belief concerning the previous states. These predictions can be used to further calculate a new possible belief, since our current belief corresponds to a distribution over \mathcal{M}_P . For each for each possible μ we determine how our beliefs would change as we acquire new observations. The main difficulty is maintaining the joint distribution over states and POMDPs.

6.3 Belief compression

Belief-augmented POMDPs in generally admit no compact representation of our current belief. This is due to the fact that there is a uncertainty both about which POMDP we are acting in and about the state of each possible POMDP. In fact, the sufficient statistic for such a problem consists of *the complete history of observations*.

This problem is not unique in reinforcement learning, however. A lot of other inference problems admit no compact posterior distributions. Gaussian processes [Rasmussen and Williams, 2006], for example, have complexity $O(n^3)$ in the number of parameters n .

For the specific case of BAPOMDPs, Ross et al [Ross et al., 2008a] give a fixed-sample approximation for the value function with a strict upper bound on the error. This work uses the idea of α -vector backups employed in Poupart et al. [2006] to evaluate the BAMDP value function on a finite POMDP. In addition to the value function approximation, the authors also employ particle filters to represent the belief. In very closely related work, Poupart and Vlassis [Poupart and Vlassis, 2008] employ sampling from reachable *beliefs* together with an efficient analytical approximation to the BAPOMDP value function.

7 Conclusion, future directions and open problems

Online methods for reinforcement learning have now reached a point of relative maturity, especially in the distribution-free framework. Methods for nearly optimal reinforcement learning now exist in the discrete case Auer et al. [2008]. The continuous case is covered only for bandit problems, however Auer et al. [2007]; Coquelin and Munos [2007]; Kocsis and Szepesvári [2006]. The continuous-case extension of the discrete framework may be relatively straightforward, but it is nevertheless unclear whether a naive extension (by a simple discretisation) of the bounds in Auer et al. [2008] will be sufficient. Related results in policy learning in continuous spaces Dimitrakakis and Lagoudakis [2008] show an exponential dependency on the number of dimensions, even though the policy iteration procedure

used therein is quite efficient. This however is probably due to the generic weakness that interval-based methods have with dealing with high dimensional spaces: they require a number of partitionings of the state space exponential in the number of dimensions. Whether such methods can be further improved remains to be seen.

There are some approaches which use sparse sampling Kearns and Singh [1998] to deal with this problem. These methods have been employed in a Bayesian online settings as well Wang et al. [2005], with promising results. Such methods, however, are unuseful when rewards are undiscounted: Kearns's sparse sampling Kearns and Singh [1998] method relies on the fact that the discount factor acts as an implicit horizon. It is possible to employ average-reward value iteration (see for example Puterman [2005]) to obtain upper bounds on the optimal value function at the leaf nodes of such a tree. However, the main problem with that approach is that the average-reward value iteration in general diverges: thus, there is no easy way to calculate the value of predecessor states.

Current research to improve the performance in the online case is mostly focused on improved methods for tree search Coquelin and Munos [2007]; Dimitrakakis [2008, 2009]; Hren and Munos [2008]; Kocsis and Szepesvári [2006]; Ross et al. [2008b]; Wang et al. [2005]. Offline methods have been explored in the context of point-based approximations to analytic solutions Poupart et al. [2006], as well as in the context of linear programming Castro and Precup [2007]. Both approaches could be effective tools to reduce computation, by allowing one to generalize over belief states. This is an approach followed by Ross et al Ross et al. [2008a] and Poupart and Vlassis [Poupart and Vlassis, 2008] for exploration in POMDPs.

In summary, the following questions should be of interest to researchers in the field: (a) How well do value function approximations on BA(PO)MDPs generalize to unseen belief states? (b) How to perform an effective discretisation of the continuous space. Can we go beyond interval-based methods? (c) How can we sample MDP and belief states efficiently? (d) Can Bayesian methods be extended to the undiscounted case?

Acknowledgments

This work was supported by the ICIS project, supported by the Dutch Ministry of Economic Affairs, grant nr: BSIK03024. Thanks to Ronald Ortner and Nikos Vlassis for interesting discussions and ideas and to the anonymous reviewers for their detailed comments and suggestions.

Appendix

Proposition 5.1 By definition, $V^*(\omega) \geq V^\pi(\omega)$ for all $\omega = (s, \xi)$, for any policy π . The lower bound follows trivially, since

$$V^{\pi^*(\bar{\mu}_\xi)}(\omega) \triangleq \int V_\mu^{\pi^*(\bar{\mu}_\xi)}(s) \xi(\mu) d\mu. \quad (7.1)$$

The upper bound is derived as follows. First note that for any function f , $\max_x \int f(x, u) du \leq \int \max_x f(x, u) du$. Then, we remark that:

$$V^*(\omega) = \max_\pi \int V_\mu^\pi(s) \xi(\mu) d\mu \quad (7.2a)$$

$$\leq \int \max_\pi V_\mu^\pi(s) \xi(\mu) d\mu \quad (7.2b)$$

$$= \int V_\mu^{\pi^*(\mu)}(s) \xi(\mu) d\mu. \quad (7.2c)$$

■

References

- Rajeev Agrawal. The continuum-armed bandit problem. *SIAM Journal on Control and Optimization*, 33(6):1926–1951, 1995. doi: 10.1137/S0363012992237273. URL <http://link.aip.org/link/?SJC/33/1926/1>.
- P. Auer, R. Ortner, and C. Szepesvari. Improved Rates for the Stochastic Continuum-Armed Bandit Problem. In *Learning Theory*, volume 4539 of *Lecture Notes in Computer Science*, page 454. Springer, 2007.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *J. Machine Learning Research*, 3(Nov):397–422, 2002.
- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2/3):235–256, 2002.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. In *Proceedings of NIPS 2008*, 2008.
- Richard Bellman and Robert Kalaba. A mathematical theory of adaptive control processes. *Proceedings of the National Academy of Sciences of the United States of America*, 45(8):1288–1290, 1959. ISSN 00278424. URL <http://www.jstor.org/stable/90152>.
- James Berger. The case for objective Bayesian analysis. *Bayesian Analysis*, 1(3):385–402, 2006.
- Dimitri Bertsekas. Dynamic programming and suboptimal control: From ADP to MPC. *Fundamental Issues in Control, European Journal of Control*, 11(4-5), 2005. From 2005 CDC, Seville, Spain.
- Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2001.
- Anselm Blumer, Andrzej Ehrenfeuch, David Haussler, and Manfred Warmuth. Learnability and the Vapnik-Chervonenkis Dimension. *Journal of the Association for Computing Machinery. Vol.*, 36(4):929–965, 1989.
- C.G.E. Boender and A.H.G. Rinnooy Kan. Bayesian stopping rules for multistart global optimization methods. *Mathematical Programming*, 37(1):59–80, 1987.
- Pablo Samuel Castro and Doina Precup. Using linear programming for bayesian exploration in Markov decision processes. In Manuela M. Veloso, editor, *IJCAI*, pages 2437–2442, 2007.
- Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning and Games*. Cambridge University press, Cambridge, UK, 2006.
- Ming-Hui Chen, Joseph G. Ibrahim, and Constantin Yiannoutsos. Prior elicitation, variable selection and Bayesian computation for logistic regression models. *Journal of the Royal Statistical Society (Series B): Statistical Methodology*, 61(1):223–242, 1999.
- Herman Chernoff. Sequential Models for Clinical Trials. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol.4*, pages 805–812. Univ. of Calif Press, 1966.

- Pierre-Arnaud Coquelin and Rémi Munos. Bandit algorithms for tree search. In *UAI '07, Proceedings of the 23rd Conference in Uncertainty in Artificial Intelligence, Vancouver, BC Canada, 2007*.
- Richard Dearden, Nir Friedman, and Stuart J. Russell. Bayesian Q-learning. In *AAAI/IAAI*, pages 761–768, 1998. URL citeseer.ist.psu.edu/dearden98bayesian.html.
- Morris H. DeGroot. *Optimal Statistical Decisions*. John Wiley & Sons, 1970.
- D. Dey, P. Muller, and D. Sinha. *Practical nonparametric and semiparametric Bayesian statistics*. Springer, 1998.
- Christos Dimitrakakis. Tree exploration for Bayesian RL exploration. In *Computational Intelligence for Modelling, Control and Automation, International Conference on*, pages 1029–1034, Wien, Austria, 2008. IEEE Computer Society. ISBN 978-0-7695-3514-2. doi: <http://doi.ieeecomputersociety.org/10.1109/CIMCA.2008.32>.
- Christos Dimitrakakis. Nearly optimal exploration-exploitation decision thresholds. In *Int. Conf. on Artificial Neural Networks (ICANN)*, 2006. IDIAP-RR 06-12.
- Christos Dimitrakakis. Complexity of stochastic branch and bound for belief tree search in Bayesian reinforcement learning. Technical Report IAS-UVA-09-01, University of Amsterdam, April 2009.
- Christos Dimitrakakis. Efficient methods for near-optimal sequential decision making under uncertainty. In Robert Babuska and Frans Groen, editors, *Interactive Collaborative Information Systems*, volume 281 of *SCI*, pages 125–153. Springer, 2010.
- Christos Dimitrakakis and Michail G. Lagoudakis. Algorithms and bounds for rollout sampling approximate policy iteration. In Girgin et al. [2008], pages 27–40. ISBN 978-3-540-89721-7.
- Christos Dimitrakakis and Christian Savu-Krohn. Cost-minimising strategies for data labelling: optimal stopping and active learning. In *Proceedings of the 5th international symposium on Foundations of Information and Knowledge Systems (FoIKS 2008)*, volume 4932 of *Lecture Notes in Computer Science*, pages 96–111, Pisa, Italy, February 2008. Springer.
- Michael O. Duff and Andrew G. Barto. Local bandit approximation for optimal learning problems. In Michael C. Mozer, Michael I. Jordan, and Thomas Petsche, editors, *Advances in Neural Information Processing Systems*, volume 9, page 1019. The MIT Press, 1997. URL citeseer.ist.psu.edu/147419.html.
- Michael O’Gordon Duff. *Optimal Learning Computational Procedures for Bayes-adaptive Markov Decision Processes*. PhD thesis, University of Massachusetts at Amherst, 2002.
- Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Action elimination and stopping conditions for the multi-armed and reinforcement learning problems. *Journal of Machine Learning Research*, pages 1079–1105, 2006.
- Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.

- Milton Friedman and Leonard J. Savage. The Utility Analysis of Choices Involving Risk. *The Journal of Political Economy*, 56(4):279, 1948.
- Milton Friedman and Leonard J. Savage. The Expected-Utility Hypothesis and the Measurability of Utility. *The Journal of Political Economy*, 60(6):463, 1952.
- Sertan Girgin, Manuel Loth, Rémi Munos, Philippe Preux, and Daniil Ryabko, editors. *Recent Advances in Reinforcement Learning, 8th European Workshop, EWRL 2008, Villeneuve d'Ascq, France, June 30 - July 3, 2008, Revised and Selected Papers*, volume 5323 of *Lecture Notes in Computer Science*, 2008. Springer. ISBN 978-3-540-89721-7.
- C. J. Gittins. *Multi-armed Bandit Allocation Indices*. John Wiley & Sons, New Jersey, US, 1989.
- Michael Goldstein. Subjective Bayesian analysis: Principles and practice. *Bayesian Analysis*, 1(3):403–420, 2006.
- Milos Hauskrecht. Value-function approximations for partially observable Markov decision processes. *Journal of Artificial Intelligence Research*, pages 33–94, Aug 2000.
- Wassily Hoeffding. Lower bounds for the expected sample size and the average risk of a sequential procedure. *The Annals of Mathematical Statistics*, 31(2):352–368, 1960. ISSN 00034851. URL <http://www.jstor.org/stable/2237951>.
- Jean-François Hren and Rémi Munos. Optimistic planning of deterministic systems. In Girgin et al. [2008], pages 151–164. ISBN 978-3-540-89721-7.
- P. Kall and S.W. Wallace. *Stochastic programming*. Wiley New York, 1994.
- Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. In *Proc. 15th International Conf. on Machine Learning*, pages 260–268. Morgan Kaufmann, San Francisco, CA, 1998. URL citeseer.ist.psu.edu/kearns98nearoptimal.html.
- F. P. Kelly. Multi-armed bandits with discount factor near one: The bernoulli case. *The Annals of Statistics*, 9(5):987–1001, September 1981.
- Levente Kocsis and Csaba Szepesvári. Bandit based Monte-Carlo planning. In *Proceedings of ECML-2006*, 2006.
- R. Duncan Luce and Howard Raiffa. *Games and Decisions*. John Wiley and Sons, 1957. Republished by Dover in 1989.
- J.J. McCall. The Economics of Information and Optimal Stopping Rules. *Journal of Business*, 38(3):300–317, 1965.
- G.V. Moustakides. Optimal stopping times for detecting changes in distributions. *Annals of Statistics*, 14(4):1379–1387, 1986.
- P. Poupart, N. Vlassis, J. Hoey, and K. Regan. An analytic solution to discrete Bayesian reinforcement learning. In *ICML 2006*, pages 697–704. ACM Press New York, NY, USA, 2006.
- Pascal Poupart and Nikos Vlassis. Model-based bayesian reinforcement learning in partially observable domains. In *International Symposium on Artificial Intelligence and Mathematics (ISAIM)*, 2008.

- Marting L. Puterman. *Markov Decision Processes : Discrete Stochastic Dynamic Programming*. John Wiley & Sons, New Jersey, US, 2005.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006. ISBN 13 978-0-262-18253-9.
- Stephane Ross, Brahim Chaib-draa, and Joelle Pineau. Bayes-adaptive POMDPs. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, Cambridge, MA, 2008a. MIT Press.
- Stéphane Ross, Joelle Pineau, Sébastien Paquet, and Brahim Chaib-draa. Online planning algorithms for POMDPs. *Journal of Artificial Intelligence Resesarch*, 32:663–704, July 2008b.
- Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In *Proc. 18th International Conf. on Machine Learning*, pages 441–448. Morgan Kaufmann, San Francisco, CA, 2001. URL citeseer.ist.psu.edu/roy01toward.html.
- Leonard J. Savage. *The Foundations of Statistics*. Dover Publications, 1972.
- T. Smith and R. Simmons. Point-based POMDP algorithms: Improved analysis and implementation. In *Proceedings of the 21st Conference on Uncertainty in Artificial Intelligence (UAI-05)*, pages 542–547, 2005.
- Robert F. Stengel. *Optimal Control and Estimation*. Dover, second edition, 1994.
- Michel Talagrand. A new look at independence. *Annals of Probability*, 24(1): 1–34, 1996.
- Vladimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 2000.
- Vladimir N. Vapnik. and Alexei Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
- Abraham Wald. *Sequential Analysis*. John Wiley & Sons, 1947. Republished by Dover in 2004.
- Tao Wang, Daniel Lizotte, Michael Bowling, and Dale Schuurmans. Bayesian sparse sampling for on-line reward optimization. In *ICML '05*, pages 956–963, New York, NY, USA, 2005. ACM. ISBN 1-59593-180-5. doi: <http://doi.acm.org/10.1145/1102351.1102472>.
- Tong Zhang. From ε -entropy to KL-entropy: Analysis of minimum information complexity density estimation. *Annals of Statistics*, 34(5):2180–2210, 2006.

Index

- action, 11
 - selection, 11
- backwards induction, 12, 15
- BAMDP, 12–14, 16
- bandit problems, 19
- BAPOMDP, 20
- Bayes' rule, 6
- Bayesian inference, 5
 - for MDPs, 12
 - for POMDPs, 19
- belief, 12, 13
 - tree, 14, 15, 18
- Bellman recursion, 12
- Bernoulli distribution, 12, 17, 19
- conjugate prior, 12, 13, 17
- control, 10, 12
- decision making, 3
 - sequential, 8
- Dirichlet distribution, 12, 13, 17
 - update, 13
- discounting, 10, 11, 18
- expectation, 3
- exploration, 19, 22
- exploration-exploitation trade-off, 2
- Hoeffding inequality, 7, 17
- horizon, 10, 11
 - finite, 14
 - infinite, 12
- hyper-state, 14
- marginal probability, 3
- Markov inequality, 7
- MDP, 11
 - unknown, 12
- mean MDP, 15–17
- Monte Carlo, 16, 17
 - belief sampling, 18, 21
 - sparse sampling, 22
 - value function lower bound, 18
 - value function upper bound, 17
- observation distribution, 19
- optimal stopping, 8, 10, 13
- policy, 11
 - blind, 16
 - greedy, 16
 - optimal, 12, 17
 - optimal for a specific MDP, 15
- POMDP, 16, 19
- preference, 4
 - action, 5
 - subjective, 4
- prior, 5
 - elicitation, 6
 - objective, 6
 - subjective, 5
- regret, 5, 7, 18, 19
 - expected, 7, 8
- reinforcement learning, 10
- reward, 11, 14, 17
 - average, 22
 - undiscounted, 18, 22
- sampling
 - experiment, 9
- state, 11
- stochastic process, 3
- transition distribution, 11, 12, 14, 16, 17, 19, 20
- uncertainty, 2–4
 - knowledge, 5
 - stochastic, 4
- upper confidence bound, 10, 18
- utility, 4, 5
 - bounds, 7
 - expected, 4
 - maximin, 4
 - risk-sensitive, 5
- value function, 11
 - bounds, 14–17
 - lower bound, 18
 - MDP, 11
 - optimal, 12
 - policy, 11