
Planning to the Information Horizon of BAMDPs via Epistemic State Abstraction

Dilip Arumugam

Department of Computer Science
Stanford University
dilip@cs.stanford.edu

Satinder Singh

DeepMind, London
baveja@deepmind.com

Abstract

The Bayes-Adaptive Markov Decision Process (BAMDP) formalism pursues the Bayes-optimal solution to the exploration-exploitation trade-off in reinforcement learning. As the computation of exact solutions to Bayesian reinforcement-learning problems is intractable, much of the literature has focused on developing suitable approximation algorithms. In this work, before diving into algorithm design, we first define, under mild structural assumptions, a complexity measure for BAMDP planning. As efficient exploration in BAMDPs hinges upon the judicious acquisition of information, our complexity measure highlights the worst-case difficulty of gathering information and exhausting epistemic uncertainty. To illustrate its significance, we establish a computationally-intractable, exact planning algorithm that takes advantage of this measure to show more efficient planning. We then conclude by introducing a specific form of state abstraction with the potential to reduce BAMDP complexity and gives rise to a computationally-tractable, approximate planning algorithm.

1 Introduction

The Bayes-Adaptive Markov Decision Process (BAMDP) [Duff, 2002] is a classic formalism encapsulating the optimal treatment of the exploration-exploitation trade-off by a reinforcement-learning agent with respect to prior beliefs over an uncertain environment. Unfortunately, the standard formulation suffers from an intractably-large hyperstate space (that is, the joint collection of environment states coupled with the agent’s current state of knowledge over the unknown environment) and much of the literature has been dedicated to identifying suitable approximations [Bellman and Kalaba, 1959, Dayan and Sejnowski, 1996, Duff and Barto, 1997, Dearden et al., 1998, Strens, 2000, Duff, 2001, 2003b,a, Wang et al., 2005, Poupart et al., 2006, Castro and Precup, 2007, Kolter and Ng, 2009, Asmuth et al., 2009, Dimitrakakis, 2009, Sorg et al., 2010, Araya-López et al., 2012, Guez et al., 2012, 2013, 2014, Ghavamzadeh et al., 2015, Zintgraf et al., 2019]. In this work, we take steps toward clarifying the hardness of BAMDPs before outlining an algorithmic concept that may help mitigate problem difficulty and facilitate near-optimal solutions.

First, we introduce the notion of *information horizon* as a complexity measure on BAMDP planning, characterizing when it is truly difficult to identify the underlying uncertain environment. Naturally, the agent’s state of knowledge at each timestep (a component of the overall BAMDP hyperstate) reflects its current epistemic uncertainty and, as the agent accumulates data, this posterior concentrates, exhausting uncertainty and identifying the true environment; after this point, the Bayes-optimal policy naturally coincides with the optimal policy of the underlying Markov Decision Process (MDP). Simply put, the information horizon quantifies the worst-case number of timesteps needed for the agent to reach this point whereupon there is no more information to be gathered about the uncertain environment.

With this complexity measure in hand, we then entertain the idea of *epistemic state abstraction* as an effective algorithmic tool for trading off between reduced information horizon (complexity) and near-Bayes-optimality of the corresponding planning solution. Intuitively, as the total number of knowledge states an agent may take on drives the intractable size of the hyperstate space, we operationalize state abstraction [Li et al., 2006, Abel et al., 2016] to perform a lossy compression of the epistemic state space, inducing a “smaller” and more tractable BAMDP for planning; our results not only mirror those of analogous work on state aggregation for improved efficiency in traditional MDP planning [Van Roy, 2006] but also parallel similar findings [Hsu et al., 2007, Zhang et al., 2012] on the effectiveness of belief state aggregation in partially-observable MDP (POMDP) planning [Kaelbling et al., 1998].

On the whole, our work provides one possible answer to a question that has already been asked and answered several times in the context of MDPs [Bartlett and Tewari, 2009, Jaksch et al., 2010, Farahmand, 2011, Maillard et al., 2014, Bellemare et al., 2016, Arumugam et al., 2021, Abel et al., 2021]: how hard is my BAMDP? While the remainder of the paper goes on to examine how one particular mechanism for reducing this complexity can translate into a more efficient planning algorithm, we anticipate that this work can serve as a starting point for building a broader taxonomy of BAMDPs, paralleling existing structural classes of MDPs [Jiang et al., 2017, Sun et al., 2019, Agarwal et al., 2020, Jin et al., 2021].

2 Problem Formulation

In this section, we formally define BAMDPs as studied in this paper. As a point of contrast, we begin by presenting the standard MDP formalism used throughout the reinforcement-learning literature [Sutton and Barto, 1998]. We use $\Delta(\mathcal{X})$ to denote the set of all probability distributions with support on an arbitrary set \mathcal{X} and denote, for any natural number $N \in \mathbb{N}$, the index set as $[N] = \{1, 2, \dots, N\}$. For any two arbitrary sets \mathcal{X} and \mathcal{Y} , we denote the class of all functions mapping from \mathcal{X} to \mathcal{Y} as $\{\mathcal{X} \rightarrow \mathcal{Y}\} \triangleq \{f \mid f : \mathcal{X} \rightarrow \mathcal{Y}\}$.

2.1 Markov Decision Processes

We begin with a sequential decision-making problem represented via the traditional finite-horizon Markov Decision Process (MDP) [Bellman, 1957, Puterman, 1994] $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \beta, H \rangle$ where \mathcal{S} is a finite set of states, \mathcal{A} is a finite set of actions, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a deterministic reward function, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is a transition function prescribing next-state transition distributions for all state-action pairs, $\beta \in \Delta(\mathcal{S})$ is an initial state distribution, and $H \in \mathbb{N}$ is the horizon denoting the agent’s total number of steps or interactions with the environment. An agent’s sequential interaction within this environment proceeds in each timestep $h \in [H]$, starting with an initial state $s_1 \sim \beta(\cdot)$, by observing the current state $s_h \in \mathcal{S}$, selecting an action $a_h \in \mathcal{A}$, and then enjoying a reward $\mathcal{R}(s_h, a_h)$ as the environment transitions to $s_{h+1} \sim \mathcal{T}(\cdot \mid s_h, a_h)$. Action selections made by the agent are governed by its non-stationary policy π : a collection of H stationary, deterministic policies $\pi = (\pi_1, \pi_2, \dots, \pi_H)$, where $\forall h \in [H]$, $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$. We quantify the performance of policy π at timestep $h \in [H]$ by its induced value function $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$ denoting the expected sum of future rewards by deploying policy

π from a particular state $s \in \mathcal{S}$: $V_h^\pi(s) = \mathbb{E} \left[\sum_{h'=h}^H \mathcal{R}(s_{h'}, a_{h'}) \mid s_h = s \right]$, where the expectation

integrates over randomness in the environment transitions. Analogously, we define the action-value function induced by policy π at timestep h as $Q_h^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ which denotes the expected future sum of rewards by being in a particular state $s \in \mathcal{S}$, executing a particular action $a \in \mathcal{A}$,

and then following policy π thereafter: $Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{h'=h}^H \mathcal{R}(s_{h'}, a_{h'}) \mid s_h = s, a_h = a \right]$. We are

guaranteed the existence of an optimal policy π^* that achieves supremal value $V_h^*(s) = \sup_{\pi \in \Pi^H} V_h^\pi(s)$

for all $s \in \mathcal{S}$, $h \in [H]$ where the policy class contains all deterministic policies $\Pi = \{\pi \mid \pi : \mathcal{S} \rightarrow \mathcal{A}\}$. Since rewards are bounded in $[0, 1]$, we have that $0 \leq V_h^\pi(s) \leq V_h^*(s) \leq H - h + 1$ for all $s \in \mathcal{S}$, $h \in [H]$, and π . These value functions obey the Bellman equation and the Bellman optimality equation, respectively:

$$V_h^\pi(s) = Q_h^\pi(s, \pi_h(s)) \quad V_h^*(s) = \max_{a \in \mathcal{A}} Q_h^*(s, a), \quad V_{H+1}^\pi(s) = 0 \quad V_{H+1}^*(s) = 0 \quad \forall s \in \mathcal{S},$$

$$Q_h^\pi(s, a) = \mathcal{R}(s, a) + \mathbb{E}_{s' \sim \mathcal{T}(\cdot | s, a)} [V_{h+1}^\pi(s')], \quad Q_h^*(s, a) = \mathcal{R}(s, a) + \mathbb{E}_{s' \sim \mathcal{T}(\cdot | s, a)} [V_{h+1}^*(s')].$$

2.2 Bayes-Adaptive Markov Decision Processes

The BAMDP formalism offers a Bayesian treatment of an agent interacting with an uncertain MDP. More specifically, a decision-making agent is faced with a MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}_\theta, \beta, H \rangle$ defined around an unknown transition function \mathcal{T}_θ ¹, for some latent parameter $\theta \in \Theta$. Prior uncertainty in θ is reflected by the distribution $p(\theta)$. In classic work on BAMDPs with finite state-action spaces, the parameters θ denote visitation counts and $p(\theta)$ is a Dirichlet distribution, so as to leverage the convenience of Dirichlet-multinomial conjugacy for exact posterior updates [Duff, 2002, Poupart et al., 2006]. For our purposes, we will assume an alternative parameterization whose importance will be made clear later when defining our complexity measure.

Assumption 1. *We assume that Θ is known and $|\Theta| < \infty$ such that an agent is only ever reasoning about its uncertainty over a finite set of $|\Theta|$ known MDPs. We further make a realizability assumption that the true parameters reside in this finite set, $\theta \in \Theta$.*

Under Assumption 1, an agent’s prior uncertainty in \mathcal{T}_θ is reflected by the distribution $p(\theta) \in \Delta(\Theta)$ which, with each step of experience encountered by the agent, may be updated via Bayes’ rule to recover a corresponding posterior distribution in light of observed data from the environment. For simplicity, we do not concern ourselves with the computation of the posterior and instead assume access to a deterministic $\mathcal{B} : \Delta(\Theta) \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \Delta(\Theta)$ that performs an exact posterior update to any input distribution $p \in \Delta(\Theta)$ based on the experience tuple $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ in $\mathcal{O}(1)$ time.

The corresponding BAMDP for \mathcal{M} is defined around a so-called *hyperstate* space $\mathcal{X} = \mathcal{S} \times \Delta(\Theta)$ such that any hyperstate $x = \langle s, p \rangle \in \mathcal{X}$ denotes the agent’s original or physical state $s \in \mathcal{S}$ within the true MDP while $p \in \Delta(\Theta)$ denotes the agent’s information state or epistemic state [Lu et al., 2021] about the uncertain environment; intuitively, the epistemic state represents the agent’s knowledge of the environment based on all previously observed data. This gives rise to the BAMDP $\langle \mathcal{X}, \mathcal{A}, \overline{\mathcal{R}}, \overline{\mathcal{T}}, \overline{\beta}, H \rangle$ where \mathcal{A} is the same action set as the original MDP \mathcal{M} , $\overline{\mathcal{R}} : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ is the same reward function as in \mathcal{M} (that is, $\overline{\mathcal{R}}(\langle s, p \rangle, a) = \mathcal{R}(s, a) \forall \langle s, p \rangle \in \mathcal{X}, a \in \mathcal{A}$), $\overline{\beta} \in \Delta(\mathcal{X})$ is defined as $\overline{\beta} = \beta \times \delta_{p(\theta)}$ where $\delta_{p(\theta)}$ denotes a Dirac delta centered around the agent’s prior $p(\theta)$, and H is the same horizon as MDP \mathcal{M} . Due to the determinism of the posterior updates given by \mathcal{B} , the BAMDP transition function $\overline{\mathcal{T}} : \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})$ is defined as

$$\overline{\mathcal{T}}(x' | x, a) = \sum_{\theta \in \Theta} \mathcal{T}_\theta(s' | s, a) p(\theta) \mathbb{1}(p' = \mathcal{B}(p, s, a, s')),$$

where $x' = \langle s', p' \rangle \in \mathcal{X}$. The associated BAMDP policy $\pi = (\pi_1, \pi_2, \dots, \pi_H)$, $\pi_h : \mathcal{X} \rightarrow \mathcal{A}, \forall h \in [H]$ selects actions based on the current state of the MDP as well as the agents accumulated knowledge of the environment thus far. With these components, we may define the associated BAMDP value functions with $x' = \langle s', \mathcal{B}(p, s, a, s') \rangle$:

$$V_h^\pi(x) = Q_h^\pi(x, \pi_h(x)) \quad V_h^*(x) = \max_{a \in \mathcal{A}} Q_h^*(x, a), \quad V_{H+1}^\pi(x) = 0 \quad V_{H+1}^*(x) = 0 \quad \forall x \in \mathcal{X},$$

$$Q_h^\pi(x, a) = \mathcal{R}(s, a) + \sum_{s', \theta} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}^\pi(x'), \quad Q_h^*(x, a) = \mathcal{R}(s, a) + \sum_{s', \theta} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}^*(x').$$

Based on these optimality equations, we see that the optimal policy of a BAMDP achieving supremal value V_h^* across all timesteps $h \in [H]$ is the Bayes-optimal policy which appropriately balances the exploration-exploitation trade-off in reinforcement learning. An observation is that this Bayes-optimal policy will tend to achieve lower value than the optimal policy of MDP \mathcal{M} as the agent takes more informative (possibly sub-optimal) actions to identify the true underlying environment.

3 Related Work

Bellman and Kalaba [1959] offer the earliest formulation of Bayesian reinforcement learning, whereby the individual actions of a decision-making agent not only provide an update to the physical state

¹For ease of exposition, we focus on uncertainty in transition dynamics although one could model uncertainty in either the reward function or the full MDP model (rewards & transitions).

of the world but also impact the agent’s internal model of how the world operates. [Dayan and Sejnowski \[1996\]](#) follow this line of thinking to derive implicit exploration bonuses based on how an agent performs posterior updates. [Kolter and Ng \[2009\]](#) make this more explicit and incorporate a specific visitation-based bonus that decays with the concentration of the agent’s Dirichlet posterior. As an alternative, [Sorg et al. \[2010\]](#) incorporate an exploration bonus based on the variance of the agent’s posterior while [Araya-López et al. \[2012\]](#) achieve optimistic exploration by boosting transition probabilities. [Duff and Barto \[1997\]](#) identify multi-armed bandits (that is, MDPs with exactly one state and arbitrarily many actions) as a unique setting where the Bayes-optimal solution is computationally tractable through the use of Gittins indices [[Gittins, 1979](#)]. While the vast space of more complicated BAMDPs are computationally intractable, a goal of this paper is to add a bit of nuance and clarify when one might still hope to recover efficient, approximate planning. This is also distinct from the PAC-BAMDP framework introduced by [Kolter and Ng \[2009\]](#), which serves as a characterization of algorithmic efficiency, rather than problem hardness.

Representing uncertainty in the optimal value function rather than environment transition function, [Dearden et al. \[1998\]](#) derive a practical Bayesian Q -learning algorithm by foregoing representation of the epistemic state and instead resampling Q^* -values at each timestep. [Strens \[2000\]](#) finds an alternate, tractable solution by lazily updating the epistemic state at the frequency of whole episodes, rather than individual timesteps; a long line of work [[Agrawal and Jia, 2017](#), [Osband et al., 2016a,b](#), [Osband and Van Roy, 2017](#), [O’Donoghue et al., 2018](#), [Osband et al., 2019](#)] analyzes this type of approximation to the Bayesian reinforcement-learning problem theoretically and also explores how to scale these solution concepts with deep neural networks.

[Duff \[2001\]](#) finds tractability in representing policies as finite-state stochastic automata, noting structural similarities between BAMDPs and partially-observable MDPs (POMDPs) [[Kaelbling et al., 1998](#)]; this type of thinking is further extended by [Poupart et al. \[2006\]](#) who exploit similar structure between the optimal value functions of BAMDPs and POMDPs. [Duff \[2003a\]](#) examine improved memory requirements when applying actor-critic algorithms [[Konda and Tsitsiklis, 2000](#)] to BAMDPs while [Duff \[2003b\]](#) consider how to approximately model the stochastic process of the evolving epistemic state via diffusion models. [Wang et al. \[2005\]](#) introduce a sparse-sampling approach [[Kearns et al., 2002](#)] for balancing computational efficiency against fidelity to Bayes-optimal action selection. An analogous sparse-sampling approach is also developed by [Castro and Precup \[2007\]](#), but with a linear-programming methodology for value-function approximation. A line of work [[Guez et al., 2012, 2013, 2014](#)] develops more scalable, sparse-sampling lookahead approaches on the back of Monte-Carlo tree search [[Kocsis and Szepesvári, 2006](#)]; these algorithms are somewhat similar in spirit to the approach of [Asmuth et al. \[2009\]](#) who merge multiple posterior samples into a single model while [Guez et al. \[2014\]](#) keep each sample distinct and integrate out the posterior randomness. For a more complete and detailed survey of Bayesian reinforcement learning, we refer readers to [Ghavamzadeh et al. \[2015\]](#). Crucially, the aforementioned approaches largely revolve around ignoring the epistemic state, lazily updating the epistemic state, or approximating the impact of the epistemic state via random sampling. In contrast, this work offers a new approach and highlights how lossy compression of the epistemic state may naturally reduce BAMDP hardness. Perhaps the most related prior work is by [Lee et al. \[2018\]](#) who introduce a practical approximate-planning approach by quantizing the epistemic state space; this paper clarifies the theoretical ramifications of this quantization step.

Our work is also connected to analyses of approximate value iteration [[Bellman, 1957](#)] in the MDP setting [[Tseng, 1990](#), [Littman et al., 1995](#)], where more recent work has managed to recover improved sample complexity bounds for approximate value iteration [[Sidford et al., 2018b,a](#)]. Like [Kearns and Singh \[1999\]](#), our algorithms utilize exact value iteration almost as a black box and it is an open question for future work to see if similar ideas and proof techniques for these approximate variants might be leveraged in the BAMDP setting. Crucially, the variants of value iteration introduced in this work are merely a backdrop for illustrating the utility of our complexity measure and, more generally, a regard for underlying information structure in BAMDPs.

The particular class of epistemic state abstraction introduced and studied in this work revolves around the covering number of the epistemic state space. Curiously, this deepens an existing connection between BAMDPs and POMDPs [[Duff, 2001](#)], where a line of work establishes the covering number of the belief state space as a viable complexity measure for the latter both in theory [[Hsu et al., 2007](#)] and in practice [[Zhang et al., 2012](#)]. In a less related but similar vein, [Kakade et al. \[2003a\]](#) establish a provably-efficient reinforcement-learning algorithm when the MDP state space is a metric space;

their corresponding sample complexity guarantee depends on the covering number of the state space under the associated metric. Our planning complexity result for abstract BAMDPs mirrors those established by these works in its dependence on the covering number of the epistemic state space.

4 The Complexity of BAMDP Planning

In this section, we examine the difficulty of solving BAMDPs through the lens of a classic planning algorithm: value iteration [Bellman, 1957]. Due to space constraints, we relegate pseudocode for all discussed algorithms to Appendix A. We begin with a quick review of the traditional algorithm applied to our setting before introducing the information horizon as a complexity measure for BAMDPs. This quantity gives rise to a more efficient planning algorithm for BAMDPs that waives excessive dependence on the original problem horizon. In order to facilitate an analysis of planning complexity in BAMDPs via value iteration, we require a finite hyperstate space \mathcal{X} . For now, we will assume that \mathcal{X} is finite, but still considerably large, by virtue of an aggressively-fine quantization of the $(|\Theta| - 1)$ -dimensional simplex, also considered in the empirical work of Lee et al. [2018]:

Assumption 2. *We assume the existence of a suitable, fixed quantization of simplex $\widehat{\Delta}(\Theta) \subset \Delta(\Theta)$ where $|\widehat{\Delta}(\Theta)| < \infty$ such that the BAMDP hyperstate space $\mathcal{X} = \mathcal{S} \times \widehat{\Delta}(\Theta)$ is finite, $|\mathcal{X}| < \infty$.*

4.1 Naive Value Iteration

To help build intuitions, we begin by presenting a typical version of value iteration for finite-horizon BAMDPs as Algorithm 1. This algorithm iterates backwards through the H timesteps, computing Q_h^* across every hyperstate-action pair. With the provision of our posterior update oracle \mathcal{B} , we avoid a square dependence on the hyperstate space ($|\mathcal{X}|^2$) and instead only require $\mathcal{O}(|\mathcal{S}||\Theta|)$ to compute next-state value. Consequently, the resulting planning complexity of Algorithm 1 is $\mathcal{O}(|\mathcal{X}||\mathcal{A}||\mathcal{S}||\Theta|H)$. Clearly, this represents an onerous burden for two distinct reasons: (1) we are forced to contend with a potentially very large horizon H and (2) we must also search through the entirety of the hyperstate space, \mathcal{X} . In the sections that follow, we alleviate the burdens of challenges (1) and (2) in series, using our new notion of information horizon to mitigate the impact of H and leveraging epistemic state abstraction to further reduce the role of $|\mathcal{X}|$, where the latter occurs at the cost of introducing approximation error.

4.2 Information Horizon

As noted in the previous section, our planning complexity suffers from its dependence on the BAMDP horizon H . A key observation, however, is that once an agent has completely resolved its uncertainty and identified one of the $|\Theta|$ environments, all that remains is to deploy the optimal policy for that particular MDP. As an exaggerated but illustrative example of this, consider a BAMDP where any action executed at the first timestep completely identifies the true environment $\theta \in \Theta$. With no residual epistemic uncertainty left, the Bayes-optimal policy would now completely coincide with the optimal policy and take actions without changing the epistemic state since, at this point, the agent has acquired all the requisite information about the previously unknown environment. Even if the problem horizon H is substantially large, a simple BAMDP like the one described should be fairly easy to solve as epistemic uncertainty is so easily diminished and information is quickly exhausted; it is this principle that underlies our hardness measure.

Let π be an arbitrary non-stationary policy. For any hyperstate $x \in \mathcal{X}$, we denote by $\mathbb{P}^\pi(x_h = x)$ the probability that policy π visits hyperstate x at timestep h . With this, we may define the reachable hyperstate space of policy π at timestep $h \in [H]$ as $\mathcal{X}_h^\pi = \{x \in \mathcal{X} \mid \mathbb{P}^\pi(x_h = x) > 0\} \subset \mathcal{X}$. In words, the reachable hyperstate space of a policy π at a particular timestep is simply the set of all possible hyperstates that may be reached by π at that timestep with non-zero probability. Recall that for any hyperstate $x = \langle s, p \rangle \in \mathcal{X}$, the epistemic state $p \in \Delta(\Theta)$ is a (discrete) probability distribution, for which we may denote its corresponding entropy as $\mathbb{H}(p)$. Given a BAMDP, we define the *information horizon of a policy* π as $\mathcal{I}(\pi) = \inf\{h \in [H] \mid \forall x_h = \langle s_h, p_h \rangle \in \mathcal{X}_h^\pi, \mathbb{H}(p_h) = 0\}$. The information horizon of a policy, if it exists, identifies the first timestep in $[H]$ where, regardless of precisely which hyperstate is reached by following π at this timestep, the agent has fully resolved all of its epistemic uncertainty over the environment θ . At this point, we call attention back to our structural Assumption 1 for BAMDPs and note that, under the standard parameterization of epistemic

state via count parameters for Dirichlet priors/posteriors, we would only be able to assess residual epistemic uncertainty through differential entropy which, unlike the traditional (Shannon) entropy $\mathbb{H}(\cdot)$, is potentially negative and has no constant lower bound [Cover and Thomas, 2012].² Naturally, to compute the *information horizon of the BAMDP*, we need only take the supremum across the non-stationary policy class: $\mathcal{I} = \sup_{\pi \in \Pi^H} \mathcal{I}(\pi)$, where $\Pi = \{\mathcal{X} \rightarrow \mathcal{A}\}$.

Clearly, when it exists, we have that $1 \leq \mathcal{I} \leq H$; the case where $\mathcal{I} = 1$ corresponds to having a prior $p(\theta)$ that is itself a Dirac delta δ_θ centered around the true environment, in which case, θ is known completely and the agent may simply compute and deploy the optimal policy for the MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}_\theta, \beta, H \rangle$. At the other end of the spectrum, an information horizon $\mathcal{I} = H$ suggests that, in the worst case, an agent may need all H steps of behavior in order to fully identify the environment. In the event that there exists any single non-stationary policy π for which the infimum of $\mathcal{I}(\pi)$ does not exist (that is, $\mathcal{I}(\pi) = \infty$), then clearly $\mathcal{I} = \infty$; this represents the most difficult, worst-case scenario wherein an agent may not always be capable of fully resolving its epistemic uncertainty within the specified problem horizon H . For certain scenarios, the supremum taken over the entire non-stationary policy class may be exceedingly strict and, certainly, creates a computational intractability should one wish to operationalize the information horizon algorithmically; in these situations, it may be more natural to consider smaller or regularized policy classes (for instance, the collection of expressible policies under a chosen neural network architecture) that yield more actionable notions of BAMDP complexity. We now go on to show how the information horizon can be used to design a more efficient BAMDP planning algorithm whose planning complexity bears a more favorable dependence on H .

4.3 Informed Value Iteration

The key insight from the previous section is that, when the information horizon exists and once an agent has acted for \mathcal{I} timesteps, the Bayes-optimal policy necessarily falls back to the optimal policy associated with the true environment. Consequently, if the solutions to all $|\Theta|$ possible underlying MDPs are computed up front, an agent can simply backup their optimal values starting from the \mathcal{I} th timestep, rather than backing up values beginning at the original horizon H . This high-level idea is implemented as Algorithm 2 which assumes access to a sub-routine `mdp_value_iteration` that consumes a MDP and produces the associated optimal value function for the initial timestep, V_1^* .

Since the underlying unknown MDP is one of $|\Theta|$ possible MDPs, Algorithm 2 proceeds by first computing the optimal value function associated with each of them in sequence using standard value iteration, incurring a time complexity of $\mathcal{O}(|\Theta| |\mathcal{S}|^2 |\mathcal{A}| (H - \mathcal{I}))$. Note that the horizon of each MDP is reduced to $H - \mathcal{I}$ acknowledging that, after identifying the true MDP in \mathcal{I} steps, an agent has only $H - \mathcal{I}$ steps of interaction remaining with the environment. With these $|\Theta|$ solutions in hand, the remainder of the algorithm proceeds with standard value iteration for BAMDPs (as in Algorithm 1), only now bootstrapping value from the \mathcal{I} timestep, rather than the original problem horizon H . Note that in Line 9, we could also compute the corresponding $\hat{\theta}$ in question by taking the mean of the next epistemic state p' , however, we use this calculation to make explicit the fact that, by definition of the information horizon, the agent has no uncertainty in θ at this point. As a result, instead of planning complexity that scales the hyperstate space size by a potentially large problem horizon, we incur a complexity of $\mathcal{O}(|\Theta| |\mathcal{S}| |\mathcal{A}| (|\mathcal{X}| \mathcal{I} + |\mathcal{S}| (H - \mathcal{I})))$. Naturally, as the gap between the information horizon \mathcal{I} and problem horizon H increases, the more favorably Algorithm 2 performs relative to the standard value iteration procedure of Algorithm 1.

In this section, we've demonstrated how the information horizon of a BAMDP has the potential to dramatically reduce the computational complexity of planning. Still, however, the corresponding guarantee bears an unfavorable dependence on the size of the hyperstate space \mathcal{X} which, in the reality that voids Assumption 2, still renders both Algorithms 1 and 2 as computationally intractable. Since this is likely inescapable for the problem of computing the exact optimal BAMDP value function, the next section considers one path for reducing this burden at the cost of only being able to realize an approximately-optimal value function.

²Prior work (see, for example, Theorem 1 of Kolter and Ng [2009]) operating with the Dirichlet parameterization will make an alternative assumption for similar effect where epistemic state updates cease after a certain number of state-action pair visitations.

5 Epistemic State Abstraction

5.1 State Abstraction in MDPs

As numerous sample-efficiency guarantees in reinforcement learning [Kearns and Singh, 2002, Kakade et al., 2003b, Strehl et al., 2009] bear a dependence on the size of the MDP state space, $|\mathcal{S}|$, a large body of work has entertained state abstraction as a tool for improving the dependence on state space size without compromising performance [Whitt, 1978, Bertsekas et al., 1988, Singh et al., 1995, Gordon, 1995, Tsitsiklis and Van Roy, 1996, Dean and Givan, 1997, Ferns et al., 2004, Jong and Stone, 2005, Li et al., 2006, Van Roy, 2006, Ferns et al., 2012, Jiang et al., 2015a, Abel et al., 2016, 2018, 2019, Dong et al., 2019, Du et al., 2019, Misra et al., 2020, Arumugam and Van Roy, 2020, Abel, 2020]. Broadly speaking, a state abstraction $\phi : \mathcal{S} \rightarrow \mathcal{S}_\phi$ maps original or ground states of the MDP into abstract states in \mathcal{S}_ϕ . Typically, one takes ϕ to be defined with respect to an abstract state space \mathcal{S}_ϕ with smaller complexity (in some sense) than \mathcal{S} ; in the case of state aggregation where all spaces in question are finite, this desideratum often takes the very simple form of $|\mathcal{S}_\phi| < |\mathcal{S}|$. Various works have identified conditions under which specific classes of state abstractions ϕ yield no approximation error and perfectly preserve the optimal policy of the original MDP [Li et al., 2006], as well as conditions under which near-optimal behavior is preserved [Van Roy, 2006, Abel et al., 2016]. As its name suggests, our proposed notion of epistemic state abstraction aims to lift these kinds of guarantees for MDPs over to BAMDPs and contend with the intractable hyperstate space.

Before examining BAMDPs, we provide a brief overview of how state abstraction impacts the traditional MDP, as a point of comparison with the BAMDP setting. Given a MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \beta, H \rangle$, a state abstraction $\phi : \mathcal{S} \rightarrow \mathcal{S}_\phi$ induces a new abstract MDP $\mathcal{M}_\phi = \langle \mathcal{S}_\phi, \mathcal{A}, \mathcal{R}_\phi, \mathcal{T}_\phi, H \rangle$ where the abstract reward function $\mathcal{R}_\phi : \mathcal{S}_\phi \times \mathcal{A} \rightarrow [0, 1]$ and transition function $\mathcal{T}_\phi : \mathcal{S}_\phi \times \mathcal{A} \rightarrow \Delta(\mathcal{S}_\phi)$ are both defined with respect to a fixed, arbitrary weighting function $\omega : \mathcal{S} \rightarrow [0, 1]$ that, intuitively, measures the contribution of each individual MDP state $s \in \mathcal{S}$ to its allocated abstract state $\phi(s)$. More specifically, ω is required to induce a probability distribution on the constituent MDP states of each abstract state: $\forall s_\phi \in \mathcal{S}_\phi, \sum_{s \in \phi^{-1}(s_\phi)} \omega(s) = 1$. This fact allows for well-defined rewards and transition probabilities as given by

$$\mathcal{R}_\phi(s_\phi, a) = \sum_{s \in \phi^{-1}(s_\phi)} \mathcal{R}(s, a)\omega(s), \quad \mathcal{T}_\phi(s'_\phi | s_\phi, a) = \sum_{s \in \phi^{-1}(s_\phi)} \sum_{s' \in \phi^{-1}(s'_\phi)} \mathcal{T}(s' | s, a)\omega(s).$$

As studied by Van Roy [2006], the weighting function ω does bear implications on the efficiency of learning and planning. Naturally, one may go on to apply various planning or reinforcement-learning algorithms to \mathcal{M}_ϕ and induce behavior in the original MDP \mathcal{M} by first applying ϕ to the current state $s \in \mathcal{S}$ and then leveraging the optimal abstract policy or abstract value function of \mathcal{M}_ϕ . Conditions under which ϕ will induce a MDP \mathcal{M}_ϕ that preserves optimal or near-optimal behavior are studied by Li et al. [2006], Van Roy [2006], Abel et al. [2016].

5.2 Compressing the Epistemic State Space

In this section, we introduce epistemic state abstraction for BAMDPs with the goal of paralleling the benefits of state abstraction in MDPs. In particular, we leverage the fact that our epistemic state space $\Delta(\Theta) = \Delta^{|\Theta|-1}$ is the $(|\Theta| - 1)$ -dimensional probability simplex. Recall that for any set \mathcal{Z} ; any threshold parameter $\delta > 0$; and any metric $\rho : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}^+$ on \mathcal{Z} , a set $\{z_1, z_2, \dots, z_K\}$ is a δ -cover of \mathcal{Z} if $\forall z \in \mathcal{Z}, \exists i \in [K]$ such that $\rho(z, z_i) \leq \delta$. In this work, we will consider δ -covers with arbitrary parameter $\delta > 0$ defined on the simplex $\Delta^{|\Theta|-1}$ with respect to the total variation distance metric on probability distributions, denoted $\|\cdot\|_{\text{TV}}$. Let $e_i \in \Delta(\Theta)$ be the i th standard basis vector such that $\mathbb{H}(e_i) = 0, \forall i \in [|\Theta|]$. We define an *epistemic state abstraction* with parameter $\delta > 0$ as the projection from $\Delta(\Theta)$ onto the smallest δ -cover of $\Delta(\Theta)$ with respect to $\|\cdot\|_{\text{TV}}$ that contains all standard basis vectors $\{e_1, e_2, \dots, e_{|\Theta|}\}$; paralleling notation for the δ -covering number, we use $\mathcal{N}(\Delta(\Theta), \delta, \|\cdot\|_{\text{TV}})$ to denote the size of this minimal cover and, for consistency with the state-abstraction literature in MDPs, use $\phi : \Delta(\Theta) \rightarrow \Delta_\phi(\Theta)$ to denote the epistemic state abstraction. Briefly, we note that while computing exact δ -covers is a NP-hard problem, approximation algorithms do exist [Hochbaum, 1996, Zhang et al., 2012]; our work here is exclusively concerned with establishing theoretical guarantees that warrant further investigation of such approximation techniques to help solve BAMDPs in practice.

It is important to note that while there are numerous statistical results expressed in terms of covering numbers (for instance, Dudley’s Theorem [Dudley, 1967]), our definition of covering number differs slightly in its inclusion of the standard basis vectors. The simple reason for this constraint is that it ensures we may still count on the existence of abstract epistemic states for which an agent has fully exhausted all epistemic uncertainty in the underlying environment. Consequently, we are guaranteed that the information horizon is still a well-defined quantity under this abstraction³. As δ increases, larger portions of the epistemic state space where the agent has residual, but still non-zero, epistemic uncertainty will be immediately mapped to the nearest standard basis vector under ϕ . If such a lossy compression is done too aggressively, the agent’s beliefs over the uncertain environment may prematurely and erroneously converge. On the other hand, if done judiciously with a prudent setting of δ , one has the potential to dramatically reduce the complexity of planning across a much smaller, *finite* hyperstate space and recover an approximately-optimal BAMDP value function.

To make this intuition more precise, consider an initial BAMDP $\langle \mathcal{X}, \mathcal{A}, \bar{\mathcal{R}}, \bar{\mathcal{T}}, \bar{\beta}, H \rangle$ and, given an epistemic state abstraction $\phi : \Delta(\Theta) \rightarrow \Delta_\phi(\Theta)$ with fixed parameter $\delta > 0$, we recover an induced abstract BAMDP $\langle \mathcal{X}_\phi, \mathcal{A}, \bar{\mathcal{R}}_\phi, \bar{\mathcal{T}}_\phi, \bar{\beta}_\phi, H \rangle$ where, most importantly, $\mathcal{X}_\phi = \mathcal{S} \times \Delta_\phi(\Theta)$ ⁴. Just as in the MDP setting, the model of the abstract BAMDP depends on a fixed, arbitrary weighting function of the original epistemic states $\omega : \Delta(\Theta) \rightarrow [0, 1]$ that adheres to the constraint: $\forall p_\phi \in \Delta_\phi(\Theta)$, $\int_{\phi^{-1}(p_\phi)} \omega(p) dp = 1$, which means abstract rewards and transition probabilities for a current and next abstract hyperstates, $x_\phi = \langle s, p_\phi \rangle$ and $x'_\phi = \langle s', p'_\phi \rangle$, are given by

$$\begin{aligned} \bar{\mathcal{R}}_\phi(x_\phi, a) &= \int_{\phi^{-1}(p_\phi)} \bar{\mathcal{R}}(x, a) \omega(p) dp = \int_{\phi^{-1}(p_\phi)} \mathcal{R}(s, a) \omega(p) dp = \mathcal{R}(s, a) \int_{\phi^{-1}(p_\phi)} \omega(p) dp = \mathcal{R}(s, a), \\ \bar{\mathcal{T}}_\phi(x'_\phi | x_\phi, a) &= \int_{\phi^{-1}(p_\phi)} \omega(p) \sum_{p' \in \phi^{-1}(p'_\phi)} \bar{\mathcal{T}}(x' | x, a) dp, \text{ where } x = \langle s, p \rangle \text{ and } x' = \langle s', p' \rangle. \end{aligned}$$

The initial abstract hyperstate distribution is defined as $\bar{\beta}_\phi = \beta \times \delta_{\phi(p(\theta))}$ where $\beta \in \Delta(\mathcal{S})$ denotes the initial state distribution of the underlying MDP while $\delta_{\phi(p(\theta))}$ is a Dirac delta centered around the agent’s original prior, $p(\theta)$, projected by ϕ into the abstract epistemic state space. Observe that the abstract BAMDP transition function is stochastic with respect to the next abstract epistemic state p'_ϕ , unlike the original BAMDP transition function whose next epistemic states are deterministic. This is, perhaps, not a surprising observation as it also occurs in standard state aggregation of deterministic MDPs as well. Nevertheless, it is important to note the corresponding abstract BAMDP value functions must now acknowledge this stochasticity for any abstract policy $\pi = (\pi_{\phi,1}, \pi_{\phi,2}, \dots, \pi_{\phi,H})$, $\pi_{\phi,h} : \mathcal{X}_\phi \rightarrow \mathcal{A}, \forall h \in [H]$:

$$\begin{aligned} V_{\phi,h}^\pi(x_\phi) &= Q_{\phi,h}^\pi(x_\phi, \pi_{\phi,h}(x_\phi)) \quad V_{\phi,h}^*(x_\phi) = \max_{a \in \mathcal{A}} Q_{\phi,h}^*(x_\phi, a), \quad V_{\phi,H+1}^\pi(x_\phi) = 0 \quad V_{\phi,H+1}^*(x_\phi) = 0 \quad \forall x_\phi \in \mathcal{X}_\phi, \\ Q_{\phi,h}^\pi(x_\phi, a) &= \mathcal{R}(s, a) + \sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(x'_\phi | x_\phi, a) V_{\phi,h+1}^\pi(x'_\phi), \quad Q_{\phi,h}^*(x_\phi, a) = \mathcal{R}(s, a) + \sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(x'_\phi | x_\phi, a) V_{\phi,h+1}^*(x'_\phi). \end{aligned}$$

Beyond the fact that this abstract BAMDP enjoys a reduced hyperstate space, we further observe that the information horizon of this new BAMDP, \mathcal{I}_ϕ , has the potential to be smaller than that of the original BAMDP. That is, if \mathcal{I} steps are needed to fully resolve epistemic uncertainty in the original BAMDP then, by compressing the epistemic state space via ϕ , we may find epistemic uncertainty exhausted in fewer than \mathcal{I} timesteps within the abstract BAMDP. Furthermore, for a suitably large setting of the δ parameter, we also have cases where the original BAMDP has $\mathcal{I} = \infty$ while $\mathcal{I}_\phi < \infty$; in words, whereas it may not have been possible to resolve all epistemic uncertainty within H timesteps, compression of the epistemic state space reduces this difficulty in the abstract problem as knowledge states near (in the total-variation sense) each vertex of the probability simplex e_i are immediately aggregated. Due to space constraints, we defer further discussion of the abstract information horizon and its relationship with the original information horizon to Appendix C.

³Note that an alternative would be to introduce an additional constant $\gamma \in \mathbb{R}^+$ and define the information horizon based on $\mathbb{H}(p) \leq \gamma$; our construction avoids carrying this cumbersome additional parameter dependence in the results.

⁴One could also imagine abstracting over the original MDP state space \mathcal{S} which, for clarity, we do not consider in this work.

As a toy illustration of last scenario, consider a ϕ with δ sufficiently large such that any step from the agent’s prior distribution immediately maps to a next abstract hyperstate with no epistemic uncertainty. Clearly, regardless of \mathcal{I} , we have an abstract BAMDP where $\mathcal{I}_\phi = 2$. Of course, under such an aggressive abstraction, we should expect to garner an unfavorable degree of approximation error between the solutions of the abstract and original BAMDPs. The next section makes this error analysis and performance loss precise alongside an approximate planning algorithm that leverages the reduced complexity of abstract BAMDPs to recover a near-optimal solution to the original BAMDP.

5.3 Informed Abstract Value Iteration

Observe that if, after inducing the abstract BAMDP according to a given epistemic state abstraction ϕ , the resulting information horizon is finite $\mathcal{I}_\phi < \infty$, then we are in a position to run Algorithm 2 on the abstract BAMDP. Moreover, we no longer need the crutch of Assumption 2 as, by definition of ϕ , we are guaranteed a finite abstract hyperstate space of size $|\mathcal{X}_\phi| = |\mathcal{S}| \cdot \mathcal{N}(\Delta(\Theta), \delta, \|\cdot\|_{\text{TV}})$. With the solution to the abstract BAMDP in hand, we can supply values to any input hyperstate of the original BAMDP $x = \langle s, p \rangle \in \mathcal{X}$ by simply applying ϕ to the agent’s current epistemic state p and querying the value of the resulting abstract hyperstate $\langle s, \phi(p) \rangle \in \mathcal{X}_\phi$. We present this approximate BAMDP planning procedure as Algorithm 3.

By construction, this algorithm inherits the planning complexity guarantee of Algorithm 2, specialized to the abstract BAMDP input, yielding $\mathcal{O}(|\Theta| |\mathcal{S}|^2 |\mathcal{A}| (\mathcal{N}(\Delta(\Theta), \delta, \|\cdot\|_{\text{TV}})^2 \mathcal{I}_\phi + (H - \mathcal{I}_\phi)))$. A key feature of this result is that we entirely forego a (direct) dependence on the hyperstate space of the original BAMDP and, instead, take on dependencies with the size of the abstract hyperstate space, $|\mathcal{X}_\phi|^2 = |\mathcal{S}|^2 \mathcal{N}(\Delta(\Theta), \delta, \|\cdot\|_{\text{TV}})^2$, and the abstract information horizon \mathcal{I}_ϕ . While both terms decrease as $\delta \rightarrow 1$, there is a delicate balance to be maintained between the ease with which one may solve the abstract BAMDP and the quality of the resulting solution when deployed in the original BAMDP of interest. We dedicate the remainder of this section to making this balance mathematically precise. Due to space constraints, all proofs are relegated to Appendix B. A natural first step in our analysis is to establish an approximation error bound:

Proposition 1. *Let V_h^* and $V_{\phi,h}^*$ denote the optimal original and abstract BAMDP value functions, respectively, for any timestep $h \in [H]$. Let ϕ be an epistemic state abstraction as defined above. Then, $\max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi,h}^*(\phi(x))| \leq 2\delta(H - h)(H - h + 1)$.*

In order to establish a complimentary performance-loss bound, we require an intermediate result characterizing performance shortfall of a BAMDP value function induced by a greedy policy with respect to another near-optimal BAMDP value function. The analogue of this result for MDPs is proven by Singh and Yee [1994], and the proof for BAMDPs follows similarly.

Proposition 2. *Let $V = \{V_1, V_2, \dots, V_H\}$ be an arbitrary BAMDP value function. We denote by $\pi_{h,V}$ the greedy policy with respect to V defined $\forall x = \langle s, p \rangle \in \mathcal{X}$ as*

$$\pi_{h,V}(x) = \arg \max_{a \in \mathcal{A}} \left[\mathcal{R}(x, a) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}(x') \right],$$

where $x' = \langle s', \mathcal{B}(p, s, a, s') \rangle \in \mathcal{X}$. Recall that V_{h+1}^* denotes the optimal BAMDP value function at timestep $h + 1$ and π_h^* denote the Bayes-optimal policy. If for all $h \in [H]$, for all $s \in \mathcal{S}$, and for any $p, q \in \Delta(\Theta)$ $|V_h^*(\langle s, p \rangle) - V_h(\langle s, q \rangle)| \leq \varepsilon$, then $\|V_h^* - V_h^{\pi_{h,V}}\|_\infty \leq 2\varepsilon(H - h + 1)$.

Combining Propositions 1 and 2 immediately yields a corresponding performance-loss bound as desired, paralleling the analogous result for state aggregation in MDPs (see Theorem 4.1 of Van Roy [2006]):

Proposition 3. *Let $\pi_{\phi,h}^*$ denote the greedy policy with respect to $V_{\phi,h+1}^*$. Then, $\|V_h^* - V_h^{\pi_{\phi,h}^*}\|_\infty \leq 4\delta(H - h)(H - h + 1)^2$.*

6 Discussion & Conclusion

In this work, we began by characterizing the complexity of a BAMDP via an upper bound on the total number of interactions needed by an agent to exhaust information and fully resolve its

epistemic uncertainty over the true environment. Under an assumption on the exact form of the agent’s uncertainty, we showed how this information horizon facilitates more efficient planning when smaller than the original problem horizon. We recognize that Assumption 1 deviates from the traditional parameterization of uncertainty in the MDP transition function via the Dirichlet distribution [Poupart et al., 2006, Kolter and Ng, 2009] (sometimes also known as the flat Dirichlet-Multinomial or FDM model [Asmuth, 2013]). The driving force behind this choice is to avoid dealing in differential entropy when engaging with the (residual) uncertainty contained in any epistemic state. Should one aspire to depart from Assumption 1 altogether in a rigorous way that manifests within the analysis, we suspect that it may be fruitful to consider a lossy compression of each epistemic state into a discrete, $|\Theta|$ -valued random variable. Under such a formulation, the appropriate tool from information theory for the analysis would be rate-distortion theory [Cover and Thomas, 2012, Csiszár, 1974]. This would, in a theoretically-sound way, allow for an arbitrary BAMDP parameterization and, for the purposes of continuing the use (discrete) Shannon entropy in the definition of the information horizon, induce a lossy compression of each epistemic state whose approximation error relative to the true epistemic state could be accounted for via the associated rate-distortion function.

Recognizing the persistence of the intractable BAMDP hyperstate space, we then proceeded to outline epistemic state abstraction as a mechanism that not only induces a finite, tractable hyperstate space but also has the potential to incur a reduced information horizon within the abstract problem. Through our analysis of approximation error and performance loss, we observe an immediate consequence of Proposition 3: if one wishes to compute an ε -optimal BAMDP value function for an original BAMDP of interest, one need only find the $\frac{\varepsilon}{4(H-h)(H-h+1)^2}$ -cover of the simplex, $\Delta(\Theta)$, and then apply the corresponding epistemic state abstraction through Algorithm 3, whose planning complexity bears no dependence on the hyperstate space of the original BAMDP and has reduced dependence on the problem horizon. One might observe that the right-hand side of the value-loss bound is maximized at timestep $h = 1$, making this first step the limiting factor when determining what value of δ to employ for computing the epistemic state abstraction. As this cover could become quite large and detract from the efficiency of utilizing an epistemic state abstraction in subsequent time periods, future work might benefit from considering abstractions formed by a sequence of exactly H δ_h -covers, where the indexing of δ_h in time $h \in [H]$ affords better preservation of value (via a straightforward extension of Proposition 3) across all timesteps simultaneously.

One caveat and limitation of our planning algorithms (both exact and approximate) is the provision of the information horizon as an input. An agent designer may seldom have the prescience of knowing the underlying BAMDP information structure or, even with suitable regularity assumptions on the policy class, be able to compute it. An observation is that many sampling-based algorithms for approximately solving BAMDPs, like BAMCP [Guez et al., 2012], implicitly hypothesize a small information horizon (typically, a value of 1) through their use of posterior sampling and choice of rollout policy. Meanwhile, recent work has demonstrated strong performance guarantees for a reinforcement-learning agent acting in an arbitrary environment [Dong et al., 2022] through the use of an incrementally increasing discount factor [Jiang et al., 2015b, Arumugam et al., 2018], gradually expanding the effective range over which the agent is expected to demonstrate competent behavior. Taking inspiration from this idea, future work might consider designing more-efficient planning algorithms that, while ignorant of the true information horizon, instead hypothesize a sequence of increasing information horizons, eventually building up to the complexity of the full BAMDP. Of course, prior to development of novel algorithms, the notion that existing BAMDP planners may already make implicit use of the information horizon is in and of itself a task for future work to tease apart and make mathematically rigorous.

Moreover, similar to how the simulation lemma [Kearns and Singh, 2002] provides a principled foundation for model-based reinforcement learning, our analysis might also be seen as offering theoretical underpinnings to the Bayes-optimal exploration strategies learned by meta reinforcement-learning agents [Ortega et al., 2019, Mikulik et al., 2020] whose practical instantiations already rely upon approximate representations of epistemic state [Zintgraf et al., 2019, 2021].

Acknowledgements

The authors gratefully acknowledge the anonymous reviewers for their insightful comments, questions, and discussions.

References

- David Abel. *A Theory of Abstraction in Reinforcement Learning*. PhD thesis, Brown University, 2020. [7](#)
- David Abel, David Hershkowitz, and Michael Littman. Near optimal behavior via approximate state abstraction. In *International Conference on Machine Learning*, pages 2915–2923. PMLR, 2016. [2](#), [7](#)
- David Abel, Dilip Arumugam, Lucas Lehnert, and Michael L. Littman. Toward good abstractions for lifelong learning. In *NeurIPS Workshop on Hierarchical Reinforcement Learning*, 2017. [22](#)
- David Abel, Dilip Arumugam, Lucas Lehnert, and Michael Littman. State abstractions for lifelong reinforcement learning. In *International Conference on Machine Learning*, pages 10–19. PMLR, 2018. [7](#)
- David Abel, Dilip Arumugam, Kavosh Asadi, Yuu Jinnai, Michael L Littman, and Lawson LS Wong. State abstraction as compression in apprenticeship learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3134–3142, 2019. [7](#)
- David Abel, Nate Umbanhowar, Khimya Khetarpal, Dilip Arumugam, Doina Precup, and Michael Littman. Value preserving state-action abstractions. In *International Conference on Artificial Intelligence and Statistics*, pages 1639–1650. PMLR, 2020. [22](#)
- David Abel, Cameron Allen, Dilip Arumugam, D. Ellis Hershkowitz, Michael L. Littman, and Lawson L.S. Wong. Bad-policy density: A measure of reinforcement learning hardness. In *ICML Workshop on Reinforcement Learning Theory*, 2021. [2](#)
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. FLAMBE: Structural complexity and representation learning of low rank MDPs. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20095–20107, 2020. [2](#)
- Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In *Advances in Neural Information Processing Systems*, pages 1184–1194, 2017. [4](#)
- Mauricio Araya-López, Vincent Thomas, and Olivier Buffet. Near-optimal BRL using optimistic local transitions. In *Proceedings of the 29th International Conference on Machine Learning*, pages 515–522, 2012. [1](#), [4](#)
- Dilip Arumugam and Benjamin Van Roy. Randomized value functions via posterior state-abstraction sampling. *arXiv preprint arXiv:2010.02383*, 2020. [7](#)
- Dilip Arumugam, David Abel, Kavosh Asadi, Nakul Gopalan, Christopher Grimm, Jun Ki Lee, Lucas Lehnert, and Michael L Littman. Mitigating planner overfitting in model-based reinforcement learning. *arXiv preprint arXiv:1812.01129*, 2018. [10](#)
- Dilip Arumugam, Peter Henderson, and Pierre-Luc Bacon. An information-theoretic perspective on credit assignment in reinforcement learning. *arXiv preprint arXiv:2103.06224*, 2021. [2](#)
- John Asmuth, Lihong Li, Michael L Littman, Ali Nouri, and David Wingate. A Bayesian sampling approach to exploration in reinforcement learning. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 19–26, 2009. [1](#), [4](#)
- John Thomas Asmuth. *Model-based Bayesian reinforcement learning with generalized priors*. Rutgers The State University of New Jersey-New Brunswick, 2013. [10](#)
- Peter L Bartlett and Ambuj Tewari. REGAL: a regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 35–42, 2009. [2](#)
- Marc G Bellemare, Georg Ostrovski, Arthur Guez, Philip Thomas, and Rémi Munos. Increasing the action gap: New operators for reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. [2](#)

- Richard Bellman. A Markovian decision process. *Journal of mathematics and mechanics*, pages 679–684, 1957. 2, 4, 5
- Richard Bellman and Robert Kalaba. On adaptive control processes. *IRE Transactions on Automatic Control*, 4(2):1–9, 1959. 1, 3
- Dimitri P Bertsekas, David A Castanon, et al. Adaptive aggregation methods for infinite horizon dynamic programming. 1988. 7
- Pablo Samuel Castro and Doina Precup. Using linear programming for Bayesian exploration in Markov decision processes. In *IJCAI*, volume 24372442, 2007. 1, 4
- Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012. 6, 10
- Imre Csiszár. On an extremum problem of information theory. *Studia Scientiarum Mathematicarum Hungarica*, 9, 1974. 10
- Peter Dayan and Terrence J Sejnowski. Exploration bonuses and dual control. *Machine Learning*, 25(1):5–22, 1996. 1, 4
- Thomas Dean and Robert Givan. Model minimization in Markov decision processes. In *AAAI/IAAI*, pages 106–111, 1997. 7
- Richard Dearden, Nir Friedman, and Stuart Russell. Bayesian Q-learning. In *Proceedings of the Fifteenth National/Tenth Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence*, pages 761–768, 1998. 1, 4
- Christos Dimitrakakis. Complexity of stochastic branch and bound methods for belief tree search in Bayesian reinforcement learning. *arXiv preprint arXiv:0912.5029*, 2009. 1
- Shi Dong, Benjamin Van Roy, and Zhengyuan Zhou. Provably efficient reinforcement learning with aggregated states. *arXiv preprint arXiv:1912.06366*, 2019. 7
- Shi Dong, Benjamin Van Roy, and Zhengyuan Zhou. Simple Agent, Complex Environment: Efficient Reinforcement Learning with Agent States. *Journal of Machine Learning Research*, 23(255):1–54, 2022. URL <http://jmlr.org/papers/v23/21-0773.html>. 10
- Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford. Provably efficient RL with rich observations via latent state decoding. In *International Conference on Machine Learning*, pages 1665–1674. PMLR, 2019. 7
- Richard M Dudley. The sizes of compact subsets of Hilbert space and continuity of Gaussian processes. *Journal of Functional Analysis*, 1(3):290–330, 1967. 8
- Michael O Duff. Monte-Carlo algorithms for the improvement of finite-state stochastic controllers: Application to Bayes-adaptive Markov decision processes. In *International Workshop on Artificial Intelligence and Statistics*, pages 93–97. PMLR, 2001. 1, 4
- Michael O Duff. Design for an optimal probe. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 131–138, 2003a. 1, 4
- Michael O Duff. Diffusion approximation for Bayesian Markov chains. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 139–146, 2003b. 1, 4
- Michael O Duff and Andrew G Barto. Local bandit approximation for optimal learning problems. In *Advances in Neural Information Processing Systems*, pages 1019–1025, 1997. 1, 4
- Michael O’Gordon Duff. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes*. University of Massachusetts Amherst, 2002. 1, 3
- Amir-massoud Farahmand. Action-gap phenomenon in reinforcement learning. *Advances in Neural Information Processing Systems*, 24:172–180, 2011. 2

- Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite Markov decision processes. In *UAI*, volume 4, pages 162–169, 2004. 7
- Norman Ferns, Pablo Samuel Castro, Doina Precup, and Prakash Panangaden. Methods for computing state similarity in Markov decision processes. *arXiv preprint arXiv:1206.6836*, 2012. 7
- Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Bayesian reinforcement learning: A survey. *Foundations and Trends® in Machine Learning*, 8(5-6):359–483, 2015. 1, 4
- John C Gittins. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(2):148–164, 1979. 4
- Geoffrey J Gordon. Stable function approximation in dynamic programming. In *Machine Learning Proceedings 1995*, pages 261–268. Elsevier, 1995. 7
- Arthur Guez, David Silver, and Peter Dayan. Efficient Bayes-adaptive reinforcement learning using sample-based search. In *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 1*, pages 1025–1033, 2012. 1, 4, 10
- Arthur Guez, David Silver, and Peter Dayan. Scalable and efficient Bayes-adaptive reinforcement learning based on monte-carlo tree search. *Journal of Artificial Intelligence Research*, 48:841–883, 2013. 1, 4
- Arthur Guez, Nicolas Heess, David Silver, and Peter Dayan. Bayes-adaptive simulation-based search with value function approximation. In *Advances in Neural Information Processing Systems*, pages 451–459, 2014. 1, 4
- Dorit S Hochbaum. Approximating covering and packing problems: set cover, vertex cover, independent set, and related problems. In *Approximation algorithms for NP-hard problems*, pages 94–143. 1996. 7
- David Hsu, Wee Sun Lee, and Nan Rong. What makes some POMDP problems easy to approximate? In *Proceedings of the 20th International Conference on Neural Information Processing Systems*, pages 689–696, 2007. 2, 4
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010. 2
- Nan Jiang, Alex Kulesza, and Satinder Singh. Abstraction selection in model-based reinforcement learning. In *International Conference on Machine Learning*, pages 179–188. PMLR, 2015a. 7
- Nan Jiang, Alex Kulesza, Satinder Singh, and Richard Lewis. The dependence of effective planning horizon on model accuracy. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 1181–1189. Citeseer, 2015b. 10
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2017. 2
- Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman Eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *arXiv preprint arXiv:2102.00815*, 2021. 2
- Nicholas K Jong and Peter Stone. State abstraction discovery from irrelevant state variables. In *IJCAI*, volume 8, pages 752–757. Citeseer, 2005. 7
- Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998. 2, 4
- Sham Kakade, Michael J Kearns, and John Langford. Exploration in metric state spaces. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 306–312, 2003a. 4
- Sham Machandranath Kakade et al. *On the sample complexity of reinforcement learning*. PhD thesis, 2003b. 7

- Michael Kearns and Satinder Singh. Finite-sample convergence rates for Q-learning and indirect algorithms. *Advances in Neural Information Processing Systems*, pages 996–1002, 1999. [4](#)
- Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002. [7](#), [10](#)
- Michael Kearns, Yishay Mansour, and Andrew Y Ng. A sparse sampling algorithm for near-optimal planning in large Markov decision processes. *Machine learning*, 49(2):193–208, 2002. [4](#)
- Levente Kocsis and Csaba Szepesvári. Bandit based Monte-Carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006. [4](#)
- J Zico Kolter and Andrew Y Ng. Near-Bayesian exploration in polynomial time. In *Proceedings of the 26th annual international conference on machine learning*, pages 513–520, 2009. [1](#), [4](#), [6](#), [10](#)
- Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in Neural Information Processing Systems*, pages 1008–1014, 2000. [4](#)
- Gilwoo Lee, Brian Hou, Aditya Mandalika, Jeongseok Lee, Sanjiban Choudhury, and Siddhartha S Srinivasa. Bayesian policy optimization for model uncertainty. In *International Conference on Learning Representations*, 2018. [4](#), [5](#)
- Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction for MDPs. *ISAIM*, 4:5, 2006. [2](#), [7](#)
- Michael L Littman, Thomas L Dean, and Leslie Pack Kaelbling. On the complexity of solving Markov decision problems. In *Proceedings of the Eleventh conference on Uncertainty in Artificial Intelligence*, pages 394–402, 1995. [4](#)
- Xiuyuan Lu, Benjamin Van Roy, Vikranth Dwaracherla, Morteza Ibrahimi, Ian Osband, and Zheng Wen. Reinforcement learning, bit by bit. *arXiv preprint arXiv:2103.04047*, 2021. [3](#)
- Odalric-Ambrym Maillard, Timothy A Mann, and Shie Mannor. How hard is my MDP?" the distribution-norm to the rescue". *Advances in Neural Information Processing Systems*, 27:1835–1843, 2014. [2](#)
- Vladimir Mikulik, Grégoire Delétang, Tom McGrath, Tim Genewein, Miljan Martic, Shane Legg, and Pedro A Ortega. Meta-trained agents implement Bayes-optimal agents. In *NeurIPS*, 2020. [10](#)
- Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. In *International conference on machine learning*, pages 6961–6971. PMLR, 2020. [7](#)
- Pedro A Ortega, Jane X Wang, Mark Rowland, Tim Genewein, Zeb Kurth-Nelson, Razvan Pascanu, Nicolas Heess, Joel Veness, Alex Pritzel, Pablo Sprechmann, et al. Meta-learning of sequential strategies. *arXiv preprint arXiv:1905.03030*, 2019. [10](#)
- Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement learning? In *International Conference on Machine Learning*, pages 2701–2710, 2017. [4](#)
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. In *Advances in neural information processing systems*, pages 4026–4034, 2016a. [4](#)
- Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. In *International Conference on Machine Learning*, pages 2377–2386, 2016b. [4](#)
- Ian Osband, Benjamin Van Roy, Daniel J Russo, and Zheng Wen. Deep exploration via randomized value functions. *Journal of Machine Learning Research*, 20(124):1–62, 2019. [4](#)
- Brendan O’Donoghue, Ian Osband, Remi Munos, and Volodymyr Mnih. The uncertainty Bellman equation and exploration. In *International Conference on Machine Learning*, pages 3836–3845, 2018. [4](#)

- Pascal Poupart, Nikos Vlassis, Jesse Hoey, and Kevin Regan. An analytic solution to discrete Bayesian reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 697–704, 2006. 1, 3, 4, 10
- Martin L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, 1994. 2
- Aaron Sidford, Mengdi Wang, Xian Wu, Lin F Yang, and Yinyu Ye. Near-optimal time and sample complexities for solving Markov decision processes with a generative model. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5192–5202, 2018a. 4
- Aaron Sidford, Mengdi Wang, Xian Wu, and Yinyu Ye. Variance reduced value iteration and faster algorithms for solving Markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 770–787. SIAM, 2018b. 4
- Satinder P Singh and Richard C Yee. An upper bound on the loss from approximate optimal-value functions. *Machine Learning*, 16(3):227–233, 1994. 9
- Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Reinforcement learning with soft state aggregation. *Advances in neural information processing systems*, pages 361–368, 1995. 7
- Jonathan Sorg, Satinder Singh, and Richard L Lewis. Variance-based rewards for approximate Bayesian reinforcement learning. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 564–571, 2010. 1, 4
- Alexander L Strehl, Lihong Li, and Michael L Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10(Nov):2413–2444, 2009. 7
- Malcolm Strens. A Bayesian framework for reinforcement learning. In *ICML*, volume 2000, pages 943–950, 2000. 1, 4
- Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on learning theory*, pages 2898–2933. PMLR, 2019. 2
- Richard S Sutton and Andrew G Barto. Introduction to reinforcement learning. 1998. 2
- Paul Tseng. Solving H-horizon, stationary Markov decision problems in time proportional to $\log(H)$. *Operations Research Letters*, 9(5):287–297, 1990. 4
- John N Tsitsiklis and Benjamin Van Roy. Feature-based methods for large scale dynamic programming. *Machine Learning*, 22(1):59–94, 1996. 7
- Benjamin Van Roy. Performance loss bounds for approximate value iteration with state aggregation. *Mathematics of Operations Research*, 31(2):234–244, 2006. 2, 7, 9
- Tao Wang, Daniel Lizotte, Michael Bowling, and Dale Schuurmans. Bayesian sparse sampling for on-line reward optimization. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 956–963, 2005. 1, 4
- Ward Whitt. Approximations of dynamic programs, I. *Mathematics of Operations Research*, 3(3): 231–243, 1978. 7
- Zongzhang Zhang, Michael Littman, and Xiaoping Chen. Covering number as a complexity measure for POMDP planning and learning. In *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012. 2, 4, 7
- Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. VariBAD: A very good method for Bayes-adaptive deep RL via meta-learning. In *International Conference on Learning Representations*, 2019. 1, 10
- Luisa M Zintgraf, Leo Feng, Cong Lu, Maximilian Igl, Kristian Hartikainen, Katja Hofmann, and Shimon Whiteson. Exploration in approximate hyper-state space for meta reinforcement learning. In *International Conference on Machine Learning*, pages 12991–13001. PMLR, 2021. 10

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Algorithms

Here we present all algorithms discussed in the main paper.

Algorithm 1 Value Iteration for BAMDPs

```

1: Input: BAMDP  $\langle \mathcal{X}, \mathcal{A}, \overline{\mathcal{R}}, \overline{\mathcal{T}}, \overline{\beta}, H \rangle$ 
2:  $V_{H+1}^*(\langle s, p \rangle) = 0, \forall \langle s, p \rangle \in \mathcal{X}$ 
3: for  $h = H, H - 1, \dots, 1$  do
4:   for  $\langle s, p \rangle \in \mathcal{X}$  do
5:     for  $a \in \mathcal{A}$  do
6:        $Q_h^*(\langle s, p \rangle, a) = \mathcal{R}(s, a) + \sum_{s', \theta} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}^*(\langle s', \mathcal{B}(p, s, a, s') \rangle)$ 
7:     end for
8:      $V_h^*(\langle s, p \rangle) = \max_{a \in \mathcal{A}} Q_h^*(\langle s, p \rangle, a)$ 
9:   end for
10: end for

```

Algorithm 2 Informed Value Iteration for BAMDPs

```

1: Input: BAMDP  $\langle \mathcal{X}, \mathcal{A}, \overline{\mathcal{R}}, \overline{\mathcal{T}}, \overline{\beta}, H \rangle$ , Information horizon  $\mathcal{I} < \infty$ 
2: for  $\theta \in \Theta$  do
3:    $V_\theta^* = \text{mdp\_value\_iteration}(\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}_\theta, \beta, H - \mathcal{I} \rangle)$ 
4: end for
5: for  $h = \mathcal{I} - 1, \mathcal{I} - 2, \dots, 1$  do
6:   for  $\langle s, p \rangle \in \mathcal{X}$  do
7:     for  $a \in \mathcal{A}$  do
8:       if  $h + 1 == \mathcal{I}$  then
9:         for  $s' \in \mathcal{S}$  do
10:           $p' = \mathcal{B}(p, s, a, s')$ 
11:           $\hat{\theta} = \sum_{\theta \in \Theta} \theta \mathbb{1}(p'(\theta) = 1)$ 
12:           $V_{h+1}^*(\langle s', p' \rangle) = V_{\hat{\theta}}^*$ 
13:        end for
14:       end if
15:        $Q_h^*(\langle s, p \rangle, a) = \mathcal{R}(s, a) + \sum_{s', \theta} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}^*(\langle s', \mathcal{B}(p, s, a, s') \rangle)$ 
16:     end for
17:      $V_h^*(\langle s, p \rangle) = \max_{a \in \mathcal{A}} Q_h^*(\langle s, p \rangle, a)$ 
18:   end for
19: end for

```

Algorithm 3 Informed Abstract Value Iteration for BAMDPs

```

1: Input: BAMDP  $\langle \mathcal{X}, \mathcal{A}, \overline{\mathcal{R}}, \overline{\mathcal{T}}, \overline{\beta}, H \rangle$ , Epistemic state abstraction  $\phi$ 
2: Induce abstract BAMDP  $\mathcal{M}_\phi = \langle \mathcal{X}_\phi, \mathcal{A}, \overline{\mathcal{R}}_\phi, \overline{\mathcal{T}}_\phi, \overline{\beta}_\phi, H \rangle$  with abstract information horizon
    $\mathcal{I}_\phi < \infty$ 
3: for  $\theta \in \Theta$  do
4:    $V_\theta^* = \text{mdp\_value\_iteration}(\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}_\theta, \beta, H - \mathcal{I}_\phi \rangle)$ 
5: end for
6: for  $h = \mathcal{I}_\phi - 1, \mathcal{I}_\phi - 2, \dots, 1$  do
7:   for  $\langle s, p_\phi \rangle \in \mathcal{X}_\phi$  do
8:     for  $a \in \mathcal{A}$  do
9:       if  $h + 1 == \mathcal{I}_\phi$  then
10:        for  $\langle s', p'_\phi \rangle \in \mathcal{X}_\phi$  do
11:           $\hat{\theta} = \sum_{\theta \in \Theta} \theta \mathbb{1}(p'_\phi(\theta) = 1)$ 
12:           $V_{\phi, h+1}^*(\langle s', p'_\phi \rangle) = V_{\hat{\theta}}^*$ 
13:        end for
14:       end if
15:        $Q_{\phi, h}^*(\langle s, p_\phi \rangle, a) = \mathcal{R}(s, a) + \sum_{s', p'_\phi} \overline{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{\phi, h+1}^*(\langle s', p'_\phi \rangle)$ 
16:     end for
17:      $V_{\phi, h}^*(\langle s, p_\phi \rangle) = \max_{a \in \mathcal{A}} Q_{\phi, h}^*(\langle s, p_\phi \rangle, a)$ 
18:   end for
19:   For any  $\langle s, p \rangle \in \mathcal{X}$ ,  $V_h^*(\langle s, p \rangle) = V_{\phi, h}^*(\langle s, \phi(p) \rangle)$ 
20: end for

```

B Proofs

B.1 Proof of Proposition 1

Proposition 1. Let V_h^* and $V_{\phi, h}^*$ denote the optimal original and abstract BAMDP value functions, respectively, for any timestep $h \in [H]$. Let ϕ be an epistemic state abstraction as defined above. Then,

$$\max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi, h}^*(\phi(x))| \leq 2\delta(H - h)(H - h + 1).$$

Proof. With a slight abuse of notation, for any hyperstate $x \in \mathcal{X}$, let $\phi(x) = \langle s, p_\phi \rangle \in \mathcal{X}_\phi$ denote its corresponding abstract hyperstate where $p_\phi = \phi(p) \in \Delta_\phi(\Theta)$. For brevity, we define $p' \triangleq \mathcal{B}(p, s, a, s')$. We have

$$\begin{aligned} \max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi, h}^*(\phi(x))| &= \max_{\langle s, p \rangle \in \mathcal{X}} \left| \max_{a \in \mathcal{A}} Q_h^*(\langle s, p \rangle, a) - \max_{a \in \mathcal{A}} Q_{\phi, h}^*(\langle s, p_\phi \rangle, a) \right| \\ &\leq \max_{\langle s, p \rangle, a \in \mathcal{X} \times \mathcal{A}} |Q_h^*(\langle s, p \rangle, a) - Q_{\phi, h}^*(\langle s, p_\phi \rangle, a)| \\ &= \max_{\langle s, p \rangle, a} \left| \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) p(\theta) V_{h+1}^*(\langle s', p' \rangle) - \sum_{s', p'_\phi} \overline{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{\phi, h+1}^*(\langle s', p'_\phi \rangle) \right| \end{aligned}$$

We now leverage the standard trick of adding “zero” by subtracting and adding the following between our two terms before applying the triangle inequality to separate them:

$$\sum_{s', p'_\phi} \overline{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{h+1}^*(\langle s', p' \rangle).$$

Examining the first term in isolation, we first observe that, by definition of the weighting function, $\int_{\phi^{-1}(p_\phi)} \omega(\overline{p}) d\overline{p} = 1$ and so we have

$$\sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) p(\theta) V_{h+1}^*(\langle s', p' \rangle) = \int_{\phi^{-1}(p_\phi)} \omega(\overline{p}) \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) p(\theta) V_{h+1}^*(\langle s', p' \rangle) d\overline{p}.$$

Expanding with the definition of the abstract BAMDP transition function, we have

$$\begin{aligned}
\sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{h+1}^*(\langle s', p' \rangle) &= \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) \bar{p}(\theta) V_{h+1}^*(\langle s', p' \rangle) \sum_{p'_\phi} \mathbb{1}(\mathcal{B}(\bar{p}, s, a, s') \in \phi^{-1}(p'_\phi)) \\
&= \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) \bar{p}(\theta) V_{h+1}^*(\langle s', p' \rangle) \sum_{p'_\phi} \mathbb{1}(\phi(\mathcal{B}(\bar{p}, s, a, s')) = p'_\phi) d\bar{p} \\
&= \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) \bar{p}(\theta) V_{h+1}^*(\langle s', p' \rangle) d\bar{p},
\end{aligned}$$

since $\phi(\mathcal{B}(\bar{p}, s, a, s'))$ belongs to exactly one abstract epistemic state. Using the fact that $V_{h+1}^*(\langle s', p' \rangle) \leq H - h$ and simplifying, we have

$$\begin{aligned}
\left| \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) p(\theta) V_{h+1}^*(\langle s', p' \rangle) - \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) \sum_{\theta, s'} \mathcal{T}_\theta(s' \mid s, a) \bar{p}(\theta) V_{h+1}^*(\langle s', p' \rangle) d\bar{p} \right| \\
\leq (H - h) \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) \sum_{\theta} |p(\theta) - \bar{p}(\theta)| d\bar{p} \\
= (H - h) \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) 2 \cdot \frac{1}{2} \sum_{\theta} |p(\theta) - \bar{p}(\theta)| d\bar{p} \\
= (H - h) \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) 2 \cdot \|p(\theta) - \bar{p}(\theta)\|_{\text{TV}} d\bar{p} \\
\leq 4\delta(H - h) \int_{\phi^{-1}(p_\phi)} \omega(\bar{p}) d\bar{p} \\
= 4\delta(H - h),
\end{aligned}$$

where the last upper bound follows from the definition of a δ -cover since $\phi(p) = \phi(\bar{p}) = p_\phi$, $\forall \bar{p} \in \phi^{-1}(p_\phi)$.

Moving on to the second term and applying Jensen's inequality, we have

$$\begin{aligned}
\left| \sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{h+1}^*(\langle s', p' \rangle) - \sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) V_{\phi, h+1}^*(\langle s', p'_\phi \rangle) \right| \\
\leq \sum_{s', p'_\phi} \bar{\mathcal{T}}_\phi(\langle s', p'_\phi \rangle \mid \langle s, p_\phi \rangle, a) \left| V_{h+1}^*(\langle s', p' \rangle) - V_{\phi, h+1}^*(\langle s', p'_\phi \rangle) \right| \\
\leq \max_{x \in \mathcal{X}} |V_{h+1}^*(x) - V_{\phi, h+1}^*(\phi(x))|.
\end{aligned}$$

Thus, putting everything together, we have established that

$$\max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi, h}^*(\phi(x))| \leq 4\delta(H - h) + \max_{x \in \mathcal{X}} |V_{h+1}^*(x) - V_{\phi, h+1}^*(\phi(x))|$$

Iterating the same sequence of steps for the latter term on the right-hand side $H - h$ more times, we arrive at a final bound

$$\max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi, h}^*(\phi(x))| \leq \sum_{\bar{h}=h}^H 4\delta(H - \bar{h}) = 4\delta \sum_{\bar{h}=1}^{H-h} \bar{h} = 4\delta \frac{(H-h)(H-h+1)}{2} = 2\delta(H-h)(H-h+1).$$

□

B.2 Proof of Proposition 2

Proposition 2. Let $V = \{V_1, V_2, \dots, V_H\}$ be an arbitrary BAMDP value function. We denote by $\pi_{h,V}$ the greedy policy with respect to V defined as

$$\pi_{h,V}(x) = \arg \max_{a \in \mathcal{A}} \left[\mathcal{R}(x, a) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}(x') \right] \quad \forall x = \langle s, p \rangle \in \mathcal{X},$$

where $x' = \langle s', \mathcal{B}(p, s, a, s') \rangle \in \mathcal{X}$. Recall that V_{h+1}^* denotes the optimal BAMDP value function at timestep $h+1$ and π_h^* denote the Bayes-optimal policy. If for all $h \in [H]$, for all $s \in \mathcal{S}$, and for any $p, q \in \Delta(\Theta)$

$$|V_h^*(\langle s, p \rangle) - V_h(\langle s, q \rangle)| \leq \varepsilon, \text{ then } \|V_h^* - V_h^{\pi_{h,V}}\|_\infty \leq 2\varepsilon(H - h + 1).$$

Proof. Fix an arbitrary timestep $h \in [H]$. For any $x \in \mathcal{X}$, define $a, \bar{a} \in \mathcal{A}$ such that $a = \pi_h^*(x)$ and $\bar{a} = \pi_{h,V}(x)$. Similarly, let $p' = \mathcal{B}(p, s, a, s')$ and $\bar{p}' = \mathcal{B}(p, s, \bar{a}, s')$. Since, by definition, $\pi_{h,V}$ is greedy with respect to V_{h+1} , we have that

$$\mathcal{R}(x, a) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, a) p(\theta) V_{h+1}(\langle s', p' \rangle) \leq \mathcal{R}(x, \bar{a}) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, \bar{a}) p(\theta) V_{h+1}(\langle s', \bar{p}' \rangle).$$

By assumption, we have that

$$V_{h+1}^*(\langle s', p' \rangle) - \varepsilon \leq V_{h+1}(\langle s', p' \rangle) \quad V_{h+1}(\langle s', \bar{p}' \rangle) \leq V_{h+1}^*(\langle s', \bar{p}' \rangle) + \varepsilon.$$

Applying both bounds to the above yields

$$\mathcal{R}(x, a) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, a) p(\theta) (V_{h+1}^*(\langle s', p' \rangle) - \varepsilon) \leq \mathcal{R}(x, \bar{a}) + \sum_{\theta, s'} \mathcal{T}_\theta(s' | s, \bar{a}) p(\theta) (V_{h+1}^*(\langle s', \bar{p}' \rangle) + \varepsilon).$$

Consequently, we have that

$$\mathcal{R}(x, a) - \mathcal{R}(x, \bar{a}) \leq 2\varepsilon + \sum_{\theta} p(\theta) \sum_{s'} V_{h+1}^*(\langle s', p' \rangle) [\mathcal{T}_\theta(s' | s, \bar{a}) - \mathcal{T}_\theta(s' | s, a)].$$

From this, it follows that

$$\begin{aligned} \|V_h^* - V_h^{\pi_{h,V}}\|_\infty &= \max_{x \in \mathcal{X}} |V_h^*(x) - V_h^{\pi_{h,V}}(x)| \\ &= \max_{x \in \mathcal{X}} |Q_h^*(x, a) - Q_h^{\pi_{h,V}}(x, \bar{a})| \\ &= \max_{x \in \mathcal{X}} \left| \mathcal{R}(x, a) - \mathcal{R}(x, \bar{a}) + \sum_{\theta} p(\theta) \sum_{s'} [\mathcal{T}_\theta(s' | s, a) V_{h+1}^*(\langle s', p' \rangle) - \mathcal{T}_\theta(s' | s, \bar{a}) V_{h+1}^{\pi_{h+1,V}}(\langle s', \bar{p}' \rangle)] \right| \\ &\leq 2\varepsilon + \max_{\langle s, p \rangle} \left| \sum_{\theta} p(\theta) \sum_{s'} \mathcal{T}_\theta(s' | s, \bar{a}) [V_{h+1}^*(\langle s', p' \rangle) - V_{h+1}^{\pi_{h+1,V}}(\langle s', \bar{p}' \rangle)] \right| \\ &\leq 2\varepsilon + \|V_{h+1}^* - V_{h+1}^{\pi_{h+1,V}}\|_\infty \\ &\vdots \\ &\leq 2\varepsilon(H - h + 1). \end{aligned}$$

where the last inequality follows by iterating the same procedure for the second term in the penultimate inequality across the remaining $H - h$ timesteps. \square

B.3 Proof of Proposition 3

Proposition 3. Let $\pi_{\phi,h}^*$ denote the greedy policy with respect to $V_{\phi,h+1}^*$. Then,

$$\|V_h^* - V_h^{\pi_{\phi,h}^*}\|_\infty \leq 4\delta(H - h)(H - h + 1)^2.$$

Proof. Since, for any $x \in \mathcal{X}$, $\phi(x)$ differs only in the epistemic state, the proof follows by realizing that the ε term of Proposition 2 is established by Proposition 1. Namely,

$$\|V_h^* - V_h^{\pi_{\phi,h}^*}\|_\infty \leq 2(H - h + 1) \max_{x \in \mathcal{X}} |V_h^*(x) - V_{\phi,h}^*(\phi(x))| \leq 4\delta(H - h)(H - h + 1)^2$$

\square

C On the Reduction of the Abstract Information Horizon

In this section, we offer two simple yet illustrative examples of BAMDPs where the use of epistemic state abstraction can either decrease or increase the information horizon of the resulting abstract BAMDP. Taken together with the main results of the paper, these examples underscore how epistemic state abstraction, similar to traditional state abstraction in MDPs, is not a panacea to sample-efficient BAMDP planning. Further investigation is needed to clarify the conditions under which an epistemic state abstraction may actually deliver upon the theoretical benefits outlined in this work.

C.1 Decreased Information Horizon under Epistemic State Abstraction

Consider a MDP whose state space is defined on the non-negative integers $\mathcal{S} = \mathbb{Z}^+ = \{0, 1, 2, \dots\}$ with two actions $\mathcal{A} = \{+, \circ\}$. For the purposes of this example, we ignore the reward function and focus on the transition function where, for any timestep $h \in [H]$, there is a fixed parameter $q > \frac{1}{2}$ such that $s_{h+1} = s_h + \Delta$ with $\Delta \sim \text{Bernoulli}(q)$ if $a_h = +$ and $\Delta \sim \text{Bernoulli}(1 - q)$ if $a_h = \circ$. In words, action $+$ has a higher probability of incrementing the agent’s current state by one whereas the \circ action is more likely to leave the agent’s state unchanged.

For the corresponding BAMDP, we take $\Theta = \{\theta_1, \theta_2\}$ where θ_1 corresponds to the true MDP transition function as described above. Meanwhile, θ_2 simply corresponds to the transition function of θ_1 with the actions flipped such that $\Delta \sim \text{Bernoulli}(q)$ if $a_h = \circ$ and $\Delta \sim \text{Bernoulli}(1 - q)$ if $a_h = +$. Clearly, when $q = 1$, epistemic uncertainty in this BAMDP is resolved immediately by the first transition whereas, with $q \downarrow \frac{1}{2}$, the two hypotheses become increasingly harder to distinguish, requiring more observed transitions from the environment and potentially exceeding the finite problem horizon H .

Now consider the epistemic state abstraction of $\Delta(\Theta)$ with parameter $\delta > 0$. Increasing $\delta \uparrow 1$ is commensurate with setting a $\gamma \in \mathbb{R}_{\geq 0}$ such that once the entropy of the current epistemic state falls below this threshold $\mathbb{H}(p_h) \leq \gamma$, we immediately have that $\phi(p_h) \in \{e_1, e_2\}$, where e_i denotes the i th standard basis vector in $\Delta(\Theta)$. By construction, any action taken in this MDP necessarily reveals information to help reduce epistemic uncertainty. Consequently, with enough observed transitions from the environment, we can use ϕ to collapse the agent’s beliefs around θ_1 with far fewer samples than what would be needed to fully exhaust epistemic uncertainty. Moreover, depending on the exact problem horizon H , this could be used to recover a finite abstract information horizon from what was an infinite information horizon in the original BAMDP.

More concretely, suppose $H = 3$ and $q = \frac{4}{5}$. While any policy will quickly stumble upon θ_1 as the most likely outcome, such a short horizon H likely does not allow the entropy of all epistemic states to fall to zero, potentially yielding an infinite information horizon $\mathcal{I} = \infty$. However, there certainly exists a value of δ such that the first two steps of behavior under any policy is sufficient for reaching a vertex state in $\Delta_\phi(\Theta)$ and identifying the underlying MDP, ultimately yielding $\mathcal{I}_\phi = 2$.

C.2 Increased Information Horizon under Epistemic State Abstraction

To show how epistemic state abstraction can work unfavorably and increase the information horizon, we consider a BAMDP where all policies rapidly resolve epistemic uncertainty, but the incorporation of an epistemic state abstraction may compromise the agent’s ability to reach a vertex of the simplex $\Delta(\Theta)$.

For any $N \in \mathbb{N}$, consider a single initial state connected to two N -state chains (an upper chain and a lower chain) with two possible actions $\mathcal{A} = \{a_1, a_2\}$. For each possible $\theta \in \Theta$, the corresponding transition function \mathcal{T}_θ will either have a_1 deterministically transition to the next state in the upper chain and have a_2 send the agent to the next state in the lower chain, or vice versa. In other words, each transition function \mathcal{T}_θ can be concisely encoded as a length N binary string $\in \{0, 1\}^N$ where, for each $n \in [N]$, the n th bit equal to 1 implies action a_1 transitions to the upper chain (immediately implying a_2 transitions to the lower chain) while a value of 0 signifies the opposite transition structure. Further suppose that, for each $n \in [N]$, the n th bit of all but one of the transition functions is identical; said differently, this structural assumptions says that, in the n th stage of the chain, there is a single informative transition that would uniquely identify the underlying MDP. A consequence of this is that, in the worst case, a BAMDP policy will only see one of the uninformative transitions in each

stage and, therefore, can only eliminate exactly one hypothesis from Θ with each timestep. Thus, for any problem horizon $H \geq N$, we are guaranteed an information horizon of N .

To see how epistemic state abstraction might inhibit planning, consider the policy that misses the informative transition in each stage. For an epistemic state abstraction with δ sufficiently large, the agent will remain stuck (via a self-looping transition) in the initial abstract epistemic state as no single uninformative transition yields sufficient information gain to move the agent through the abstract epistemic state space. This phenomenon of state abstraction ameliorating generalization while drastically worsening the challenge of exploration has already been observed in the standard MDP setting [Abel et al., 2017, 2020]; here, we see that epistemic state abstraction is also vulnerable to the same weakness. As a result, the agent will never converge to one of the vertices of the simplex and never see its epistemic uncertainty in the underlying environment completely diminished, resulting in $\mathcal{I}_\phi = \infty$.