

---

# On learning history-based policies for controlling Markov decision processes

---

**Gandharv Patil**  
McGill University, Mila

**Aditya Mahajan**  
McGill University, Mila

**Doina Precup**  
McGill University, Mila

## Abstract

Reinforcement learning (RL) folklore suggests that methods of function approximation based on history, such as recurrent neural networks or state abstractions that include past information, outperform those without memory, because function approximation in Markov decision processes (MDP) can lead to a scenario akin to dealing with a partially observable MDP (POMDP). However, formal analysis of history-based algorithms has been limited, with most existing frameworks concentrating on features without historical context. In this paper, we introduce a theoretical framework to examine the behaviour of RL algorithms that control an MDP using feature abstraction mappings based on historical data. Additionally, we leverage this framework to develop a practical RL algorithm and assess its performance across various continuous control tasks.

## 1 Introduction

State abstraction and function approximation are vital components used by reinforcement learning (RL) algorithms to efficiently solve complex control problems when exact computations are intractable due to the size of the state and action spaces. Over the past few decades, state abstraction in RL has evolved from the use of pre-determined and problem-specific features Crites and Barto [1995], Tsitsiklis and Roy [1996], Bertsekas and Tsitsiklis [1996], Sutton and Barto [1998], Singh et al. [2002], Kwok and Fox [2004], Proper and Tadepalli [2006] to adaptive basis functions learnt by solving an isolated regression problem [Ormonet and Sen, 2002,

Menache et al., 2005, Keller et al., 2006, Petrik, 2007], and more recently to neural network-based Deep-RL algorithms that compute state embeddings using successive layers of a neural network [Barto et al., 2004, Bellemare et al., 2019].

There has also been a renewed interest in using history-based feature abstraction methods, as growing empirical evidence suggests that such methods are beneficial in practice [OpenAI et al., 2019]. However, a theoretical characterisation of history-based Deep-RL algorithms for fully observable Markov Decision Processes (MDPs) is largely absent from the literature.

In this paper, we bridge this gap between theory and practice by providing a theoretical analysis of history-based RL agents acting in a MDP.

Our approach adapts the notion of approximate information state (AIS) for Partially Observable Markov Decision Processes (POMDPs) proposed in Subramanian et al. [2020], Subramanian and Mahajan [2019] to feature abstraction in MDPs. We also develop a theoretically grounded policy search algorithm for history-based feature abstractions and policies, which we evaluate empirically.

The rest of the paper is organised as follows: In Section 2, following a brief review of feature-based abstraction, we motivate the need for using history-based feature abstractions. In Section 3, we present a AIS-based model for the co-design of the feature abstraction and control policy, and derive its corresponding dynamic programming decomposition. We also derive bounds on the quality of approximate solutions to this dynamic program. In Section 4 we build on these approximation bounds to develop an RL algorithm for learning a history-based state representation and control policy. In Section 5, we present an empirical evaluation of our proposed algorithm on continuous control tasks. Finally, we discuss related work in Section 6 and conclude with future research directions in Section 7.

## 2 Background and Motivation

Consider an MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, r, \gamma \rangle$  where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  denotes the action space,  $P$  denotes the controlled transition matrix,  $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  denotes the per-step reward, and  $\gamma \in (0, 1)$  denotes the discount factor.

The performance of a randomised (and possibly history-dependent) policy  $\pi$  starting from a start state  $s_0$  is measured by the value function, defined as:

$$V^\pi(s_0) = \mathbb{E}^\pi \left[ \sum_{t=1}^{\infty} \gamma^{t-1} r(S_t, A_t) \mid S_0 = s_0 \right]. \quad (1)$$

A policy maximising  $V^\pi(s_0)$  over all (randomised and possibly history dependent) policies is called an optimal policy with respect to initial state  $s_0$  and is denoted by  $\pi^*$ .

In many applications,  $\mathcal{S}$  and  $\mathcal{A}$  are combinatorially large or uncountable, which makes it intractable to compute the optimal policy.

Most practical RL algorithms overcome this hurdle by using function approximation where the state is mapped to a feature space  $\mathcal{Z}$  using a state abstraction function  $\phi: \mathcal{S} \rightarrow \mathcal{Z}$ . In Deep-RL algorithms, the last layer of the network is often viewed as a feature vector. These feature vectors are then used as an approximate state for approximating the value function  $\tilde{V}: \mathcal{Z} \rightarrow \mathcal{R}$  and/or computing an approximately optimal policy  $\mu: \mathcal{Z} \rightarrow \Omega(\mathcal{A})$  [Sutton and Barto, 1998] (where  $\Omega(\mathcal{A})$  denotes the set of probability distribution over actions). Therefore, the mapping from state to distribution of actions is given by the ‘‘flattened’’ policy  $\tilde{\mu} = \mu \circ \phi$  *i.e.*,  $\tilde{\mu} = \mu(\phi(\cdot))$ .

A well known fact about function approximation is that the features that are used as an approximate state may not satisfy the controlled Markov property *i.e.*, in general,

$$\mathbb{P}(Z_{t+1} \mid Z_{1:t}, A_{1:t}) \neq \mathbb{P}(Z_{t+1} \mid Z_t, A_t).$$

To see the implications of this fact, consider the toy MDP depicted in Figures 1a to 1c, with  $\mathcal{S} = \{0, 1, 2, 3\}$ ,  $\mathcal{A} = \{0, 1, 2\}$ ,  $\{P_{s,s'}(a)\}_{a \in \mathcal{A}}$ , and  $r(0) = r(1) = -1$ ,  $r(2) = 1$ ,  $r(3) = -K$ , where  $K$  is a large positive number. Given the reward structure the objective of the policy is to try to avoid state 3 and keep the agent at state 2 as much as possible. It is easy to see that the optimal policy is

$$\pi^*(0) = 0, \quad \pi^*(1) = 0, \quad \pi^*(2) = 1, \quad \text{and} \quad \pi^*(3) = 2.$$

Note that if the initial state is not state 3 then an agent will never visit that state under the optimal

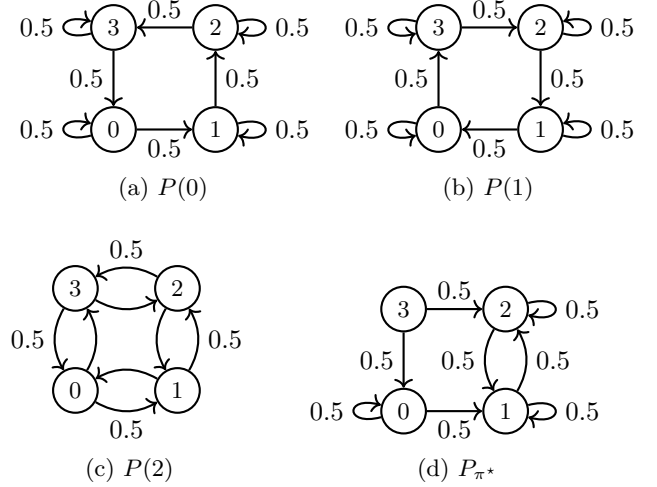


Figure 1: The transition probability for an example MDP

policy. Furthermore, any policy which cannot prevent the agent from visiting state 3 will have a large negative value and, therefore, cannot be optimal.

Now suppose the feature space  $\mathcal{Z} = \{0, 1\}$ . It is easy to see that for any Markovian feature-abstraction  $\phi: \mathcal{S} \rightarrow \mathcal{Z}$ , no policy  $\hat{\pi}: \mathcal{Z} \rightarrow \mathcal{A}$  can prevent the agent from visiting state 3. Thus, the best policy when using Markovian feature abstraction will perform significantly worse than the optimal policy (which has direct access to the state).

However, it is possible to construct a history-based feature-abstraction  $\phi$  and a history-based control policy  $\hat{\pi}$  that works with  $\phi$  and is of the same quality as  $\pi^*$ . For this, consider the following *codebooks* (where the entries denoted by a dot do not matter):

Now define

$$D(0) = \begin{bmatrix} 0 & 1 \\ 1 & 2 \\ 2 & 3 \\ 3 & 0 \end{bmatrix}, \quad F(0) = \begin{bmatrix} 0 & 1 & \cdot & \cdot \\ \cdot & 0 & 1 & \cdot \\ \cdot & \cdot & 0 & 1 \\ 1 & \cdot & \cdot & 0 \end{bmatrix},$$

$$D(1) = \begin{bmatrix} 3 & 0 \\ 0 & 1 \\ 1 & 2 \\ 2 & 3 \end{bmatrix}, \quad F(1) = \begin{bmatrix} 1 & \cdot & \cdot & 0 \\ 0 & 1 & \cdot & \cdot \\ \cdot & 0 & 1 & \cdot \\ \cdot & \cdot & 0 & 1 \end{bmatrix},$$

$$D(2) = \begin{bmatrix} 1 & 3 \\ 0 & 2 \\ 1 & 3 \\ 0 & 2 \end{bmatrix}, \quad F(2) = \begin{bmatrix} \cdot & 0 & \cdot & 1 \\ 0 & \cdot & 1 & \cdot \\ \cdot & 0 & \cdot & 1 \\ 0 & \cdot & 1 & \cdot \end{bmatrix},$$

and consider the feature-abstraction policy  $Z_t = F_{S_{t-1}, S_t}(A_{t-1})$  and a control policy  $\mu$  which is a finite state machine with memory, where the memory

$M_t$  that is updated as  $M_t = D_{M_{t-1}, Z_t}(A_{t-1})$  and the action  $A_t$  is chosen as  $A_t = \pi(M_t)$ , where  $\pi: \mathcal{S} \rightarrow \Omega(\mathcal{A})$  is any pre-specified reference policy. It can be verified that if the system starts from a known initial state then  $\mu \circ \phi = \pi$ . Thus, if we choose the reference policy  $\pi = \pi^*$ , then the agent will never visit state 3 under  $\mu \circ \phi$ , in contrast to Markovian feature-abstraction policies where (as we argued before) state 3 is always visited.

In the above example, we used the properties of the system dynamics and the reward function to design a history-based feature abstraction which outperforms memoryless feature abstractions. We are interested in developing such history-based feature abstractions using a learning framework when the system model is not known. We present such a construction in the next section.

### 3 Approximation bounds for history-based feature abstraction

The approximation results in this work depend on the properties of metrics on probability spaces. Therefore, we will first provide a brief overview of a general class of metrics known as Integral Probability Metrics (IPMs) [Müller, 1997]. Note that many commonly used metrics on probability spaces, such as total variation (TV) distance, Wasserstein distance, and maximum-mean discrepancy (MMD), are instances of IPMs.

#### 3.1 Integral Probability Metrics (IPM)

**Definition 3.1** ([Müller, 1997]). Let  $(\mathcal{E}, \mathcal{G})$  be a measurable space and let  $\mathfrak{F}$  denote a class of uniformly bounded measurable functions on  $(\mathcal{E}, \mathcal{G})$  such that the set  $\mathfrak{F}$  is convex and balanced. The Integral Probability Metric (IPM) between two probability distributions  $\nu_1, \nu_2 \in \mathcal{P}(\mathcal{E})$  with respect to the function class  $\mathfrak{F}$  is defined as:

$$d_{\mathfrak{F}}(\nu_1, \nu_2) \triangleq \sup_{f \in \mathfrak{F}} \left| \int_{\mathcal{E}} f d\nu_1 - \int_{\mathcal{E}} f d\nu_2 \right|. \quad (2)$$

For any function  $f$  (not necessarily in  $\mathfrak{F}$ ), the Minkowski functional  $\rho_{\mathfrak{F}}$  associated with the metric  $d_{\mathfrak{F}}$  is defined as:

$$\rho_{\mathfrak{F}}(f) \triangleq \inf\{\rho \in \mathcal{R}_{\geq 0} : \rho^{-1}f \in \mathfrak{F}\}. \quad (3)$$

Eq. (3), implies that for any function  $f$ :

$$\left| \int_{\mathcal{E}} f d\nu_1 - \int_{\mathcal{E}} f d\nu_2 \right| \leq \rho_{\mathfrak{F}}(f) d_{\mathfrak{F}}(\nu_1, \nu_2). \quad (4)$$

In this paper, we use the following IPMs:

1. **Total Variation Distance:** The IPM is Total Variation Distance when  $\mathfrak{F} = \mathfrak{F}^{\text{TV}} \triangleq \{f : \|f\|_{\infty} \leq 1\}$ . The Minkowski functional associated with  $d_{\mathfrak{F}}$ , denoted as  $\rho_{\mathfrak{F}^{\text{TV}}}(f)$  given as  $\rho_{\mathfrak{F}^{\text{TV}}}(f) = \frac{1}{2} \text{span}(f)$ , where  $\frac{1}{2} \text{span}(f) \triangleq \frac{1}{2}(\max(f) - \min(f))$ .
2. **Wasserstein/Kantorovich-Rubinstein Distance:** The IPM is the Wasserstein or the Kantorovich distance, when  $\mathfrak{F} = \mathfrak{F}^W \triangleq \{f : L_f \leq 1\}$  (where  $L_f$  denotes the Lipschitz constant of  $f$  with respect to the metric on some metric space  $\mathcal{E}$ ). The Minkowski function for the Wasserstein distance is  $\rho_{\mathfrak{F}^W}(f) = L_f$ .
3. **Maximum Mean Discrepancy (MMD) Distance:** The IPM is a Maximum Mean Discrepancy distance when  $\mathfrak{F}$  is  $\mathfrak{F}^{\text{MMD}} \triangleq \{f \in \mathcal{U} : \|f\|_{\mathcal{U}} \leq 1\}$  where  $\mathcal{U}$  is a reproducing kernel Hilbert space (RKHS) of real-valued functions on  $\mathcal{E}$ , and  $\|\cdot\|_{\mathcal{U}}$  denotes the associated RKHS norm. The Minkowski functional is associated with MMD is  $\rho_{\mathfrak{F}^{\text{MMD}}}(f) = \|f\|_{\mathcal{U}}$ .

#### 3.2 Approximate information state

Given an MDP  $\mathcal{M}$  and a feature space  $\mathcal{Z}$ , let  $\mathcal{H}_t = \mathcal{S} \times \mathcal{A}$  denote the space of all histories  $(S_{1:t}, A_{1:t-1})$  up to time  $t$ . Here  $S_{1:t}$  is a concise notation for the history of states  $(S_1, \dots, S_t)$ , and a similar interpretation applies to  $A_{1:t}$ . Our objective is to learn history-based feature abstraction functions  $\{\sigma_t : \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$  along with a time homogeneous policy  $\mu : \mathcal{Z} \rightarrow \Omega(\mathcal{A})$  such that the flattened policy  $\pi = \{\pi_t\}_{t \geq 1}$ , where  $\pi_t = \mu \circ \sigma_t$ , is approximately optimal.

**Definition 3.2.** A family of history-based feature abstraction functions  $\{\sigma_t : \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$  are said to be *recursively updatable* if there exists an update function  $\hat{f} : \mathcal{Z} \times \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{Z}$  such that the process  $\{Z_t\}_{t \geq 1}$ , where  $Z_t = \sigma_t(S_{1:t}, A_{1:t-1})$ , satisfies the following:

$$Z_{t+1} = \hat{f}(Z_t, S_{t+1}, A_t). \quad t \geq 1 \quad (5)$$

Since the feature abstraction approximates the state, its quality depends on how well it can be used to approximate the per step reward and predict the next state. The definition below formalises this intuition.

**Definition 3.3.** Given a family of recursively updatable history based feature abstraction functions  $\{\sigma_t : \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$ , the features  $Z_t = \sigma_t(S_{1:t}, A_{1:t-1})$  are said to be  $(\epsilon, \delta)$ -*approximate information state* (AIS) with respect to a function space  $\mathfrak{F}$  if there exist: (i) a reward approximation function  $\hat{r} : \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{R}$ , and (ii) an approximate transition kernel  $\hat{P} : \mathcal{Z} \times \mathcal{A} \rightarrow \Omega(\mathcal{S})$  such that  $Z$  satisfies the following properties:

(P1) Sufficient for approximate performance evaluation: for all  $t$ ,

$$|r(S_t, A_t) - \hat{r}(Z_t, A_t)| \leq \epsilon. \quad (6)$$

(P2) Sufficient for predicting future states approximately: for all  $t$

$$d_{\mathfrak{F}}(P(\cdot|S_t, A_t), \hat{P}(\cdot|Z_t, A_t)) \leq \delta. \quad (7)$$

We call the tuple  $(\hat{r}, \hat{P})$  as an  $(\epsilon, \delta)$ - AIS approximator. Note that similar definitions have appeared in other works *e.g.*, latent state [Gelada et al., 2019], and approximate information state for POMDPs [Subramanian et al., 2020, Subramanian and Mahajan, 2019]. However, in [Gelada et al., 2019] it is assumed that the feature abstractions are memoryless and the discussion is restricted to Wasserstein distance.

The **main difference** from the POMDP model in [Subramanian et al., 2020, Subramanian and Mahajan, 2019] is that, in POMDPs the observation  $Z_t$  is a pre-specified function of the state while in the proposed model  $Z_t$  depends on our choice of feature abstraction.

As such, the key insight presented in this paper is that an AIS-approximator of a recursively updatable history-based feature abstraction can be used to define a dynamic program. In particular, given a history-based abstraction function  $\{\sigma_t : \mathcal{H}_t \rightarrow \mathcal{Z}\}_{t \geq 1}$  which is recursively updatable using  $\hat{f}$  and an  $(\epsilon, \delta)$  AIS-approximator  $(\hat{P}, \hat{r})$ , we can define the following dynamic programming decomposition:

For any  $z_t \in \mathcal{Z}$ ,  $a_t \in \mathcal{A}$

$$\begin{aligned} \hat{Q}(z_t, a_t) &= \hat{r}(z_t, a_t) \\ &+ \gamma \sum_{s_{t+1} \in \mathcal{S}} \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(z_t, s_{t+1}, a_t)) \end{aligned} \quad (8)$$

$$\hat{V}(z_t) = \max_{a_t \in \mathcal{A}} \hat{Q}(z_t, a_t), \quad \forall z_t \in \mathcal{Z} \quad (9)$$

**Definition 3.4.** Define  $\mu : \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  be any policy such that for any  $z \in \mathcal{Z}$ ,

$$\text{Supp}(\mu(z)) \subseteq \arg \max_{a \in \mathcal{A}} \hat{Q}(z, a). \quad (10)$$

Since  $\mu$  is a policy that maps from the feature space to actions, we can use it to define a policy that maps from the history of state action pairs to actions as:

$$\pi_t(s_{1:t}, a_{1:t-1}) \triangleq \mu(\sigma_t(s_{1:t}, a_{1:t-1})) \quad (11)$$

Therefore, the dynamic program defined in (9) indirectly defines a history-based policy  $\pi = (\pi_1, \pi_2, \dots)$ . The performance of any such history-based policy is given by the following dynamic program:

For any  $h_t \in \mathcal{H}$ ,  $a_t \in \mathcal{A}$

$$\begin{aligned} Q_t^\pi(h_t, a_t) &= r(s_t, a_t) + \\ &\gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V_{t+1}^\pi(h_{t+1}), \end{aligned} \quad (12)$$

$$V_t^\pi(h_t) = \sum_{a_t \in \mathcal{A}} \pi(a_t|h_t) Q_t^\pi(h_t, a_t), \quad \forall h_t \in \mathcal{H} \quad (13)$$

We are interested in quantifying the loss in performance when such a history-based policy is used to control the MDP. Note that since  $V_t^\pi$  is not time-homogeneous, we need to compute the worst-case difference between  $V^*$  and  $V_t^\pi$ , which is given by:

$$\Delta \triangleq \sup_{t \geq 0} \sup_{h_t = (s_{1:t}, a_{1:t}) \in \mathcal{H}_t} |V^*(s_t) - V_t^\pi(h_t)|, \quad (14)$$

Our main approximation result is the following:

**Theorem 3.5.** *The worst case difference between  $V^*$  and  $V_t^\pi$  is bounded by*

$$\Delta \leq 2 \frac{\epsilon + \gamma \delta \kappa_{\mathfrak{F}}(\hat{V}, \hat{f})}{1 - \gamma}, \quad (15)$$

where  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f}) = \sup_{z, a} \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z, a)))$ ,  $\rho_{\mathfrak{F}}(\cdot)$  is the Minkowski functional associated with the IPM  $d_{\mathfrak{F}}$  as defined in (3).

Proof in Appendix A

Some salient features of the bound are as follows: First, the bound depends on the choice of metric on probability spaces. Different IPMs will result in a different value of  $\delta$  and also a different value of  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f})$ . Second, the bound depends on the properties of  $\hat{V}$ . For this reason we call it an instance dependent bound. Sometimes, it is desirable to have bounds which do not require solving the dynamic program in (9). We present such bounds as below, note that these ‘‘instance independent’’ bounds are the derived by upper bounding  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f})$ . Therefore, these are looser than the upper bound in Theorem 3.5

**Corollary 3.6.** *If the function class  $\mathfrak{F}$  is  $\mathfrak{F}^{TV}$ , then  $\Delta$  as defined in (14) is upper bounded as:*

$$\Delta \leq \frac{2\epsilon}{(1 - \gamma)} + \frac{\gamma \delta \text{span}(\hat{r})}{(1 - \gamma)^2}. \quad (16)$$

Proof in Appendix B

**Corollary 3.7.** *Let  $L_{\hat{r}}$  and  $L_{\hat{P}}$  denote the Lipschitz constants of the approximate reward function  $\hat{r}$  and approximate transition function  $\hat{P}$  respectively, and  $L_{\hat{f}}$  is the uniform bound on the Lipschitz constant of  $\hat{f}$*

<sup>1</sup>We have dropped the time index from the policy to reduce clutter

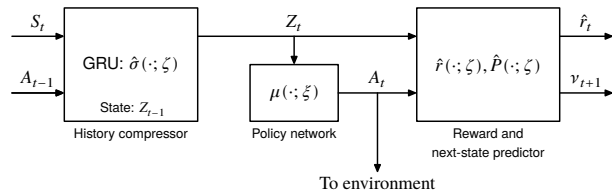


Figure 2: AIS approximator block

with respect to the state  $S_t$ . If  $\gamma L_{\hat{P}} L_{\hat{f}} \leq 1$  and the function class  $\mathfrak{F}$  is  $\mathfrak{F}^W$ , then  $\Delta$  as defined in (14) is upper bounded as:

$$\Delta \leq \frac{2\epsilon}{(1-\gamma)} + \frac{2\gamma\delta L_{\hat{r}}}{(1-\gamma)(1-\gamma L_{\hat{f}} L_{\hat{P}})}. \quad (17)$$

Proof in Appendix C

**Corollary 3.8.** If the function class  $\mathfrak{F}$  is  $\mathfrak{F}^{MMD}$ , then  $\Delta$  as defined in (14) is upper bounded as:

$$\Delta \leq 2 \frac{\epsilon + \gamma\delta\kappa_{\mathcal{U}}(\hat{V}, \hat{f})}{(1-\gamma)}, \quad (18)$$

where  $\mathcal{U}$  is a RKHS space,  $\|\cdot\|_{\mathcal{U}}$  its associated norm and  $\kappa_{\mathcal{U}}(\hat{V}, \hat{f}) = \sup_{z,a} \|\hat{V}(\hat{f}(\cdot, z, a))\|_{\mathcal{U}}$ .

*Proof.* The proof follows from the properties of MMD described previously.  $\square$

In the following section we will show how one can use these theoretical insights to design a policy search algorithm.

## 4 Reinforcement learning with history-based feature abstraction

In this section, we will use the ideas presented in the previous section to develop a computational framework that utilizes a reinforcement learning (RL) algorithm to simultaneously learn an Abstract Information State (AIS) and a policy.

The previous section helps us establish two results. The first result indicates that a history-based representation can be considered an AIS if it is capable of evolving like a state and can approximately predict the instantaneous reward and state transition. The second result suggests that if a policy is derived using an AIS, then its performance loss is bounded by the approximation error as in (15). To ensure that the implementation meets the definition of the AIS, the main idea is to represent both the AIS generator and the policy using a parametric family of functions/distributions and train them using a multi-timescale optimization algorithm.

According to Definition 3.3, the AIS generator consists of four components: a compression function  $\sigma_t$ , the update function  $\hat{f}$ , an approximate reward predictor  $\hat{r}$ , and a transition kernel  $\hat{P}$ . We can represent the history compression function using any time series approximators, such as LSTMs or GRUs. An advantage of such memory-based neural networks is that their internal layers are updated in a state-like manner. Therefore, we can satisfy Definition 3.2 since  $Z_t$  evolves according to the RNN’s state update function such that  $\hat{f} : \mathcal{Z} \times \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{Z}$ .

The main function of  $\hat{r}$  and  $\hat{P}$  is to ensure that  $Z_t$  satisfies properties P1 and P2 (in Definition 3.3), *i.e.*, prediction of the instantaneous reward with a bounded error  $\epsilon$  and approximation of the ground MDP’s transition function with a bounded error  $\delta$ . One way in which the computational framework can satisfy these conditions is by explicitly optimising the AIS generator for the constants  $\epsilon$  and  $\delta$ . We can achieve this by modelling the reward predictor  $\hat{r}$  using a multilayered perceptron (MLP) layer which uses the representation  $Z_t$  and action  $A_t$  to predict the reward  $\hat{r}_t$ . In the same way, we can model the approximate transition kernel  $\hat{P}$  using an appropriate class of stochastic kernel approximators *e.g.*, a softmax function or a mixture of Gaussian’s to learn a parametric approximation of  $P$ . We can then train the AIS generator by minimising an appropriate objective function.

To make things more concrete, let us denote the AIS generator as the following collection:  $\{\sigma_t(\cdot; \zeta), \hat{f}(\cdot; \zeta), f_{\hat{r}}(\cdot; \zeta), f_{\hat{P}}(\cdot; \zeta)\}$  where  $f_{\hat{r}}(\cdot; \zeta) : \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{R}$  and  $f_{\hat{P}}(\cdot; \zeta) : \mathcal{Z} \times \mathcal{A} \rightarrow \Omega(\mathcal{S})$  are the reward and transition approximators, and  $\zeta$  are the parameters of the respective sub-components. Instead of optimising the reward prediction loss and the transition loss separately, we can combine them in a single objective function objective function as follows:

$$\mathcal{L}_{\text{AIS}}(\zeta) = \frac{1}{T} \sum_{t=0}^T \left( \lambda \underbrace{(\hat{r}(Z_t, A_t; \zeta) - r(S_t, A_t))^2}_{\mathcal{L}_{\hat{R}}(\cdot; \zeta)} + (1-\lambda) \cdot \underbrace{d_{\mathfrak{F}}(\hat{P}(Z_t, A_t; \zeta), P)^2}_{\mathcal{L}_{\hat{P}}(\cdot; \zeta)} \right), \quad (19)$$

where,  $T$  is the length of the episode or the rollout length,  $\lambda \in [0, 1]$  is a hyper-parameter, reward prediction loss  $\mathcal{L}_{\hat{R}}(\cdot; \zeta)$  is simply the mean-squared error between the predicted and the observed reward, whereas the transition prediction loss  $\mathcal{L}_{\hat{P}}(\cdot; \zeta)$  is the distance between predicted and observed transition distributions  $\hat{P}$  and  $P$ . To compute  $\mathcal{L}_{\hat{P}}(\cdot; \zeta)$ , we need to choose an IPM. In principle we can pick any IPM, but we would want to use an IPM using which the distance

$d_{\mathfrak{F}}$  can be efficiently computed.

#### 4.1 Choice of an IPM

To compute the IPM  $d_{\mathfrak{F}}$  we need to know the probability density functions  $\hat{P}$  and  $P$ . As we assume that  $\hat{P}$  belongs to a parametric family, we know its density function in closed form. However, since we are in the learning setup, we can only access samples from  $P$ . For a function  $f \in \mathfrak{F}$ , and probability density functions  $P$  and  $\hat{P}$  such that,  $\nu_1 = P$ , and  $\nu_2 = \hat{P}$ , we can estimate the IPM  $d_{\mathfrak{F}}$  between a distribution and samples using the duality  $|\int_{\mathcal{Z}} f d\nu_1 - \int_{\mathcal{Z}} f d\nu_2|$ . In this paper, we use two forms of IPMs, the MMD distance and the Wasserstein/Kantorovich–Rubinstein distance.

##### 4.1.1 MMD Distance:

Let  $m_{\zeta}$  denote the mean of the distribution  $\hat{P}(\cdot; \zeta)$ . Then, the AIS-loss when MMD is used as an IPM is given by

$$\mathcal{L}_{\text{AIS}}(\zeta) = \frac{1}{T} \sum_{t=0}^T \left( \lambda (\hat{r}(Z_t, A_t; \zeta) - r(S_t, A_t))^2 + (1 - \lambda) (m_{\zeta}^{S_t} - 2S_t)^{\top} m_{\zeta}^{S_t} \right), \quad (20)$$

where  $m_{\zeta}^{S_t}$  is obtained using the from the transition approximator, *i.e.*, the mapping  $\hat{P}(\zeta) : \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{R}$ . For the detailed derivation of the above loss see Appendix D.1.1

##### 4.1.2 Wasserstein/Kantorovich–Rubinstein distance:

In principle, the Wasserstein/Kantorovich distance can be computed by solving a linear program [Sriperumbudur et al., 2012], but doing at every episode can be computationally expensive.

Therefore, we propose to approximate the Wasserstein distance using a KL-divergence [Kullback and Leibler, 1951] based upper-bound. The simplified-KL divergence based AIS loss is given as:

$$\mathcal{L}_{\text{AIS}}(\zeta) = \frac{1}{T} \sum_{t=0}^T \left( \lambda (\hat{r}(Z_t, A_t; \zeta) - r(S_t, A_t))^2 + (1 - \lambda) \log(\hat{P}(S_t; \zeta)) \right), \quad (21)$$

where after dropping the terms which do not depend on  $\zeta$ , we get  $d_{\mathfrak{F}}^2(P, \hat{P}) \leq \log(\hat{P}(S_t; \zeta))$  is the simplified-KL- divergence based upper bound. For the details of this derivation see Appendix D.1.2.

## 4.2 Policy gradient algorithm

---

### Algorithm 1: Policy Search with AIS

---

**Input** :  $\iota_0$ : Initial state distribution,  
 $\zeta_0$ : Ais parameters,  
 $\xi_0$ : Actor parameters,  
 $a_0$ : Initial action,  
 $\mathcal{D} = \emptyset$ : Replay buffer,  
 $N_{\text{comp}}$ : Computation budget,  
 $N_{\text{ep}}$ : Episode length,  
 $N_{\text{grad}}$ : Gradient steps

```

1 for iterations  $i = 0 : N_{\text{comp}}$  do
2   Sample start state  $s_0 \sim \iota_0$ ;
3   for iterations  $j = 0 : N_{\text{ep}}$  do
4      $z_j = \sigma_{\zeta}(s_{1:j}, a_{1:j-1})$ ;
5      $a_j = \mu_{\xi}(z_j)$ ;
6      $s_{j+1} = P(s_j, a_j)$ ;
7      $\mathcal{D} \leftarrow \{z_j, a_j, s_j, s_{j+1}\}$ ;
8      $a_{j-1} = a_j$ ;
9      $s_j = s_{j+1}$ ;
10  end
11  for every batch  $b \in \mathcal{D}$  do
12    for gradient step  $t = 0 : N_{\text{grad}}$  do
13       $\zeta_{t+1,b} = \zeta_{t,b} + b \nabla_{\zeta} \mathcal{L}_{\text{AIS}}(\zeta_{t,b})$ ;
14       $\xi_{t+1,b} = \xi_{t,b} + d \hat{\nabla}_{\xi} J(\xi_{t,b}, \zeta_{t,b})$ 
15    end
16  end
17 end

```

---

Following the design of the AIS block, we now provide a policy- gradient algorithm for learning both the AIS and the policy. The schematic of our agent architecture is given in Figure 2, and pseudo-code is given in Algorithm 1. Given a feature space  $\mathcal{Z}$ , we can simultaneously learn the AIS- generator and the policy using a multi-timescale stochastic gradient ascent algorithm [Borkar, 2008]. Let  $\mu(\cdot; \xi) : \mathcal{Z} \rightarrow \Delta(\mathcal{A})$  be a parameterised stochastic policy with parameters  $\xi$ . Let  $J(\xi, \zeta)$  denote the performance of the policy  $\mu(\cdot; \xi)$ . The policy gradient theorem [Sutton et al., 1999, Williams, 2004, Baxter and Bartlett, 2001] states that:

For a rollout horizon  $T$ , we can estimate  $\nabla_{\xi} J$  as:

$$\hat{\nabla}_{\xi} J(\xi_t, \zeta_t) = \sum_{t=1}^T \gamma^{t-1} r_t \left( \sum_{\tau=1}^t \nabla_{\xi} \log(\mu(A_t | Z_t; \xi_t)) \right).$$

Following a rollout of length  $T$ , we can then update the parameters  $\{(\zeta_i, \xi_i)\}_{i \geq 1}$  as follows:

$$\zeta_{i+1} = \zeta_i + b_i \nabla_{\zeta} \mathcal{L}_{\text{AIS}}(\zeta_i), \quad (22)$$

$$\xi_{i+1} = \xi_i + d_i \hat{\nabla}_{\xi} J(\xi_i, \zeta_i) \quad (23)$$

where the step-size  $\{b_i\}_{i \geq 0}$  and  $\{d_i\}_{i \geq 0}$  satisfy the standard conditions  $\sum_i b_i = \infty$ ,  $\sum_i b_i^2 < \infty$ ,  $\sum_i d_i =$

$\infty$  and  $\sum_i d_i^2 < \infty$  respectively. Moreover, one can ensure that the AIS generator converges faster by choosing an appropriate learning rates such that,  $\lim_{i \rightarrow \infty} \frac{d_i}{b_i} = 0$ .

### 4.3 Actor Critic Algorithm

We can also use the aforementioned ideas to design an AIS based actor-critic algorithm. In addition to a parameterised policy  $\pi(\cdot; \xi)$  and AIS generator  $(\sigma_t(\cdot; \zeta), \hat{f}, \hat{r}, \hat{P})$  the actor-critic algorithm uses a parameterised critic  $\hat{V}(\cdot; \vartheta) : \mathcal{Z} \rightarrow \mathcal{R}$ , where  $\vartheta$  are the parameters for the critic. The performance of policy  $\mu(\cdot; \xi)$  is then given by  $J(\xi, \zeta, \vartheta)$ . According to policy gradient theorem [Sutton et al., 1999, Baxter and Bartlett, 2001] the gradient of  $J(\xi, \zeta, \vartheta)$ , is given as:

$$\nabla_{\xi} J(\xi, \zeta, \vartheta) = \mathbb{E} \left[ \nabla_{\xi} \log(\mu(\cdot; \xi)) \hat{V}(\cdot; \vartheta) \right]. \quad (24)$$

And for a trajectory of length  $T$ , we approximate it as:

$$\hat{\nabla}_{\xi} J(\xi, \zeta, \vartheta) = \frac{1}{T} \sum_{t=1}^T \left[ \nabla_{\xi} \log(\mu(\cdot; \xi)) \hat{V}(\cdot; \vartheta) \right]. \quad (25)$$

The parameters  $\vartheta$  can be learnt by optimising the temporal difference loss given as:

$$\mathcal{L}_{\text{TD}}(\xi, \zeta, \vartheta) = \frac{1}{T} \sum_{t=0}^T \text{smoothL1}(\hat{V}(Z_t; \vartheta) \quad (26)$$

$$- r(Z_t, A_t) - \gamma \hat{V}(Z_{t+1}; \vartheta)). \quad (27)$$

The parameters  $\{(\zeta_i, \xi_i, \vartheta_i)\}_{i \geq 1}$  can be updated using a multi-timescale stochastic approximation algorithm as follows:

$$\zeta_{i+1} = \zeta_i + b_i \nabla_{\zeta} \mathcal{L}_{\text{AIS}}(\zeta_i), \quad (28a)$$

$$\vartheta_{i+1} = \vartheta_i + c_i \nabla_{\vartheta} \mathcal{L}_{\text{TD}}(\xi_i, \zeta_i, \vartheta_i), \quad (28b)$$

$$\xi_{i+1} = \xi_i + d_i \hat{\nabla}_{\xi} J(\xi_i, \zeta_i, \vartheta), \quad (28c)$$

where the step-size  $\{b_i\}_{i \geq 0}$ ,  $\{c_i\}_{i \geq 0}$  and  $\{d_i\}_{i \geq 0}$  satisfy the standard conditions  $\sum_i b_i = \infty$ ,  $\sum_i b_i^2 < \infty$ ,  $\sum_i c_i = \infty$ ,  $\sum_i c_i^2 < \infty$ ,  $\sum_i d_i = \infty$  and  $\sum_i d_i^2 < \infty$  respectively. Moreover, one can ensure that the AIS generator converges first, followed by the critic and the actor by choosing an appropriate step-sizes such that,  $\lim_{i \rightarrow \infty} \frac{d_i}{b_i} = 0$  and  $\lim_{i \rightarrow \infty} \frac{c_i}{d_i} = 0$ .

We elaborate more on convergence of both the methods in Appendix E.

## 5 Empirical evaluation

Through our experiments, we seek to answer the following questions:

- (1) Can history-based feature representations policies help improve the quality of solution found by memory-less RL algorithms?
- (2) In terms of the solution quality how does the proposed method compare with other methods which use memory augmented policies as well as reward and transition predictors?
- (3) How does the choice of IPM affect the algorithms performance?

We answer question (1) by comparing our approach with the proximal policy gradient (PPO) [Schulman et al., 2017] and the policy-gradient version of DeepMDP framework [Gelada et al., 2019]. For question (2) we compare our approach with modified versions of PlaNet [Hafner et al., 2019], Dreamer [Hafner et al., 2020].

For question (3) we compare the performance of our method using different MMD kernels and KL-divergence based approximation of Wasserstein distance (included in Appendices F.4 and F.5). All the approaches are evaluated on six continuous control tasks from the MuJoCo [Todorov et al., 2012] OpenAI-Gym suite. To ensure a fair comparison, the baselines and their respective hyper-parameter settings are taken from well tested stand-alone implementations provided by Dhariwal et al. [2017].

From an implementation perspective, our framework can be used to modify any off-the-shelf policy-gradient algorithm by simply replacing (or augmenting) the feature abstraction layers of the policy and/or value networks with recurrent neural networks (RNNs), trained with the appropriate losses, as outlined previously. In these experiments, we replace the fully connected layers in PPO’s architecture with a Gated Recurrent Unit (GRU). For all the implementations, we initialise the hidden state of the GRU to zero at the beginning of the trajectory. This strategy simplifies the implementation and also allows for independent decorrelated sampling of sequences, therefore ensuring robust optimisation of the networks [Hausknecht and Stone, 2015].

The results are summarised in Figure 3 (mean return with interquartile range (shaded region) estimated over 50 independent runs). The proposed algorithm consistently learns faster and achieves higher final performance than the baselines. We note that as reported in the Appendix F, MMD is a better kernel to use.

It is important to note that we can extend our framework to other policy gradient methods such as SAC [Haarnoja et al., 2018], TD3 [Fujimoto et al., 2018] or DDPG [Lillicrap et al., 2016], after satisfying certain technical conditions. However, we leave these extensions for future work. Additional experimental details and results on different IPMs can be found in Appendix F.

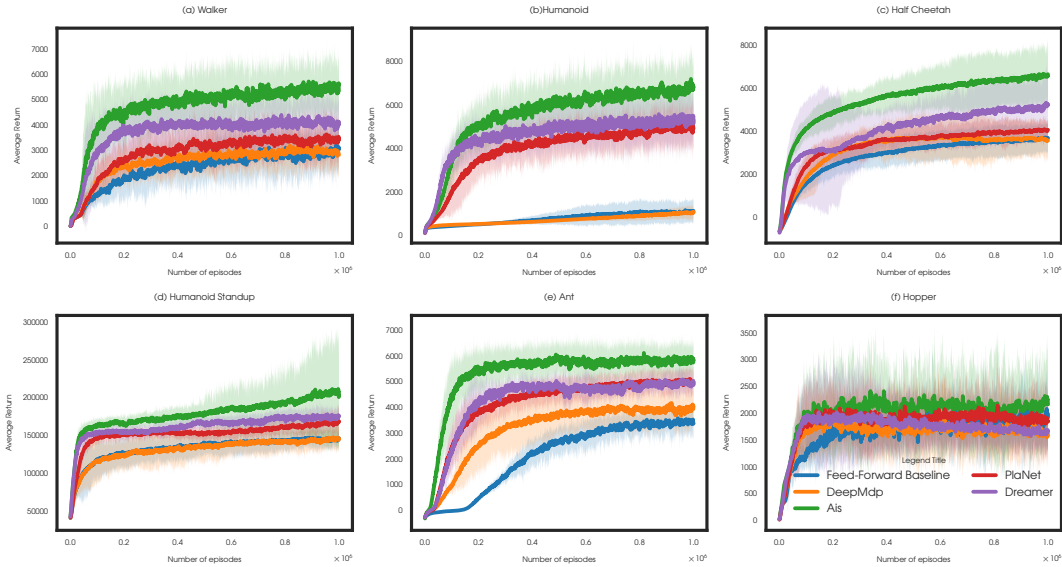


Figure 3: Empirical results averaged over 50 Monte Carlo runs with shaded regions representing the quantiles.

## 6 Related Work

The development of RL algorithms with memory-based feature abstractions has been a active area of research. Most existing algorithms have approached this challenge using non-parametric methods, such as Nearest Neighbor [Bentley, 1975, Friedman et al., 1977, Peng, 1995], Locally-weighted Regression [Baird and Klopff, 1993, Atkeson et al., 1997, Moore et al., 1997], and Kernel-based Regression [Connell and Utgoff, 1987, Dietterich and Wang, 2001, Ormonet and Sen, 2002, Xu et al., 2006, Bhat et al., 2012, Barreto et al., 2016]. Despite their robust theoretical foundations, these methods have limited applicability due to their scalability challenges in high-dimensional state and action spaces. More recently, several methods employing recurrent neural networks (RNNs) for learning history-based abstractions have demonstrated significant success in complex computer games [Hausknecht and Stone, 2015, Jaderberg et al., 2017, Espeholt et al., 2018, Gruslys et al., 2018, Ha and Schmidhuber, 2018]. However, most of these methods were designed for partially observable environments where the use of history-based methods is often essential. To the best of our knowledge, the only other work that applies a history-based RL algorithm for controlling a Markov Decision Process (MDP) is presented by OpenAI et al. [2019], where the authors illustrate that an LSTM-based agent architecture yields superior performance in object reorientation tasks using robotic arms. Nevertheless, they do not offer a theoretical analysis of their approach. In a recent contribution, Bolland et al. [2023] introduce optimisation by continuation, demonstrating that the incorporation of historical information can also enhance the

convergence properties of policy optimisation methods. Establishing concrete connection between optimisation by continuation framework and the ideas in this paper is an interesting future research direction.

### 6.1 Bisimulation metrics

On the theoretical front, our work is closely related to state aggregation techniques based on bisimulation metrics, as proposed by Givan et al. [2003], Ferns et al. [2004, 2011]. The bisimulation metric represents the fixed point of an operator within the space of semi-metrics over the state space of an MDP, characterised by Lipschitz continuity in value functions. Beyond state aggregation, bisimulation metrics have found applications in feature discovery [Comanici and Precup, 2011, Ruan et al., 2015] and transfer learning [Castro and Precup, 2010]. Despite their potential, computational challenges have hindered widespread adoption. Our contribution can be seen as offering an alternative to bisimulation for analysing history-based state abstractions and deep RL methods. Furthermore, our work extends the DeepMDP framework [Gelada et al., 2019], adapting it for history-based policies and direct policy search methods.

### 6.2 AIS and Agent state

The notion of AIS is closely related to the epistemic state recently proposed by Lu et al. [2021]. An epistemic state is a bounded representation of the history. It is updated recursively as the agent collects more information, and is represented as an environment proxy  $\Upsilon$  which is learnt by optimising a target/objective func-



tion  $\chi$ . Since  $\Upsilon$  is a random variable, its entropy  $\mathbb{H}(\Upsilon)$  is used to represent system’s uncertainty about the environment. The framework proposed in this paper can be considered as a practical way of constructing the system epistemic state where, the AIS  $Z_t$  represents both the epistemic state and the environment proxy  $\Upsilon$ ,  $\mathcal{L}_{\text{AIS}}$  represents  $\chi$ , and instead of entropy, the constants  $\epsilon$ , and  $\delta$  represent the system’s uncertainty about the environment. Studying the AIS framework in the regret minimisation paradigm is an interesting research direction as it can help establish concrete links between  $\epsilon$ ,  $\delta$ , and  $\mathbb{H}(\Upsilon)$ . This exploration has the potential to guide designers in developing principled algorithms that integrate concepts such as information-directed sampling into direct policy search algorithms.

### 6.3 Analysis of RL algorithms with attention mechanism

Recently, there has been considerable interest in developing RL algorithms which use attention mechanism/transformer architectures [Bahdanau et al., 2015, Xu et al., 2015] for learning feature abstractions [Zambaldi et al., 2019, Mott et al., 2019, Sorokin et al., 2015, Oh and Kaneko, 2018, Ritter et al., 2021, Parisotto et al., 2020, Chen et al., 2021, Loynd et al., 2020, Tang et al., 2020, Pritzel et al., 2017]. Attention mechanism extracts task relevant information from historical observations and can be used instead of RNNs for processing sequential data [Vaswani et al., 2017]. As we do not impose a functional form on the history compression function  $\sigma_t(\cdot)$  in Definition 3.3, any attention mechanism can be interpreted as history compression function, and one can construct a valid information state by ensuring that the output of the attention mechanism satisfies (P1) and (P2). That being said, even without optimising  $\mathcal{L}_{\text{AIS}}$ , the approximation bound in Theorem 3.5 still applies for RL algorithms with attention mechanisms, with the caveat that the constants  $\epsilon$ , and  $\delta$  may be arbitrarily large. A thorough empirical analysis of the effect of different attention mechanisms, and the AIS loss on the error constants  $\epsilon$ , and  $\delta$  could help us gain a better understanding of the way in which such design choices could influence the learning process.

### 6.4 AIS for POMDPs

The concept of an AIS used in this paper is similar to the idea of AIS for POMDPs [Subramanian and Mahajan, 2019, Subramanian et al., 2020]. Moreover, the literature also contains several other methods which have enjoyed empirical success in using history-based policies for controlling POMDPs [Holmes and Jr., 2006, Daswani et al., 2013, Hutter, 2014, Schaefer et al., 2007, Hafner et al., 2020, 2019, Lambrecht et al., 2023].

In principle, one can use any of these methods for controlling MDPs. However, this does not immediately provide a tight bound for the approximation error. The MDP model has more structure than POMDPs, and our goal in this paper is to use this fact to present a tighter analysis of the approximation error.

## 7 Conclusion and future work

We developed and analysed a principled approach for learning history-based policies for controlling MDPs. Such policies are useful when using function approximation, as it introduces partial observability. We believe that our approximation bounds can be helpful for practitioners to study the effect of some of their design choices on the solution quality. On the practical side, the proposed algorithm shows favourable results on high-dimensional control tasks. Note that one can also use the bounds in Theorem 3.5 to analyse the approximation error of other history-based methods. However, since some of these algorithms do not satisfy Definition 3.3, the resulting approximation error might be arbitrarily large. Such blow-ups in the approximation error could be due to the bound itself being loose or to the optimality gap being large. A sharper analysis of the approximation error by considering the specific design choices of other methods is an interesting direction for future research. Another interesting direction would be to conduct a further empirical evaluation exploring the design choices of history compression functions.

## Acknowledgements

We would like to thank Jayakumar Subramanian for the initial conceptualisation of this idea and to the anonymous reviewers for helpful comments. This research was enabled in part by support provided by Calcul Quebec and the Digital Research Alliance of Canada.

## References

- C. G. Atkeson, A. W. Moore, and S. Schaal. *Locally Weighted Learning*, pages 11–73. Springer Netherlands, Dordrecht, 1997. ISBN 978-94-017-2053-3. doi: 10.1007/978-94-017-2053-3\_2. URL [https://doi.org/10.1007/978-94-017-2053-3\\_2](https://doi.org/10.1007/978-94-017-2053-3_2).
- D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.0473>.

- L. C. Baird and A. H. Klopff. Reinforcement learning with high-dimensional, continuous actions. Technical report, Wright Laboratory, 1993.
- A. Barreto, P. Doina, and P. Joelle. Practical kernel-based reinforcement learning. *Journal of Machine Learning Research*, 2016.
- A. Barto, C. Anderson, and R. Sutton. Synthesis of nonlinear control surfaces by a layered associative search network. *Biological Cybernetics*, 43:175–185, 2004.
- J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. *J. Artif. Intell. Res.*, 15:319–350, 2001. doi: 10.1613/jair.806. URL <https://doi.org/10.1613/jair.806>.
- M. G. Bellemare, W. Dabney, R. Dadashi, A. A. Taïga, P. S. Castro, N. L. Roux, D. Schuurmans, T. Lattimore, and C. Lyle. A geometric perspective on optimal representations for reinforcement learning. In H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 4360–4371, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/3cf2559725a9fdfa602ec8c887440f32-Abstract.html>.
- J. L. Bentley. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18: 509–517, 1975.
- D. P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1st edition, 1996. ISBN 1886529108.
- N. Bhat, C. C. Moallemi, and V. F. Farias. Non-parametric approximate dynamic programming via the kernel method. In *NIPS*, 2012.
- A. Bolland, G. Louppe, and D. Ernst. Policy gradient algorithms implicitly optimize by continuation. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=3Ba6Hd3nZt>.
- V. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge University Press, 2008. ISBN 9780521515924. URL <https://books.google.ca/books?id=QLxIvgAACAAJ>.
- V. S. Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29:291–294, 1997.
- G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- P. S. Castro and D. Precup. Using bisimulation for policy transfer in mdps. In *AAAI*, 2010.
- L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *CoRR*, abs/2106.01345, 2021. URL <https://arxiv.org/abs/2106.01345>.
- G. Comanici and D. Precup. Basis function discovery using spectral clustering and bisimulation metrics. In *AAAI*, 2011.
- M. E. Connell and P. Utgoff. Learning to control a dynamic physical system. *Computational Intelligence*, 3, 1987.
- R. H. Crites and A. G. Barto. Improving elevator performance using reinforcement learning. In D. S. Touretzky, M. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8, NIPS, Denver, CO, USA, November 27-30, 1995*, pages 1017–1023. MIT Press, 1995. URL <https://proceedings.neurips.cc/paper/1995/file/390e982518a50e280d8e2b535462ec1f-Paper.pdf>.
- M. Daswani, P. Sunehag, and M. Hutter. Q-learning for history-based reinforcement learning. In C. S. Ong and T. B. Ho, editors, *Asian Conference on Machine Learning, ACML 2013, Canberra, ACT, Australia, November 13-15, 2013*, volume 29 of *JMLR Workshop and Conference Proceedings*, pages 213–228. JMLR.org, 2013. URL <http://proceedings.mlr.press/v29/Daswani13.html>.
- P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, Y. Wu, and P. Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- T. G. Dietterich and X. Wang. Batch value function approximation via support vectors. In *NIPS*, 2001.
- L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu. IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1407–1416. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/espeholt18a.html>.

- N. Ferns, P. Panangaden, and D. Precup. Metrics for finite markov decision processes. In *Conference on Uncertainty in Artificial Intelligence*, 2004.
- N. Ferns, P. Panangaden, and D. Precup. Bisimulation metrics for continuous markov decision processes. *SIAM J. Comput.*, 40:1662–1714, 2011.
- J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3:209–226, 1977.
- S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. In J. G. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1582–1591. PMLR, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>.
- K. Fukumizu, A. Gretton, G. Lanckriet, B. Schölkopf, and B. K. Sriperumbudur. Kernel choice and classifiability for rkhs embeddings of probability distributions. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009. URL <https://proceedings.neurips.cc/paper/2009/file/685ac8cad1be5ac98da9556bc1c8d9e-Paper.pdf>.
- C. Gelada, S. Kumar, J. Buckman, O. Nachum, and M. G. Bellemare. Deepmdp: Learning continuous latent space models for representation learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2170–2179. PMLR, 2019. URL <http://proceedings.mlr.press/v97/gelada19a.html>.
- R. Givan, T. L. Dean, and M. Greig. Equivalence notions and model minimization in markov decision processes. *Artif. Intell.*, 147:163–223, 2003.
- A. Gruslys, W. Dabney, M. G. Azar, B. Piot, M. G. Bellemare, and R. Munos. The reactor: A fast and sample-efficient actor-critic agent for reinforcement learning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=rkHVZWAZ>.
- D. Ha and J. Schmidhuber. World models. *CoRR*, abs/1803.10122, 2018. URL <http://arxiv.org/abs/1803.10122>.
- T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In J. G. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1856–1865. PMLR, 2018. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- D. Hafner, T. P. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent dynamics for planning from pixels. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2555–2565. PMLR, 2019. URL <http://proceedings.mlr.press/v97/hafner19a.html>.
- D. Hafner, T. P. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=S110TC4tDS>.
- M. J. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In *2015 AAAI Fall Symposia, Arlington, Virginia, USA, November 12-14, 2015*, pages 29–37. AAAI Press, 2015. URL <http://www.aaai.org/ocs/index.php/FSS/FSS15/paper/view/11673>.
- M. P. Holmes and C. L. I. Jr. Looping suffix tree-based inference of partially observable hidden state. In W. W. Cohen and A. W. Moore, editors, *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*, volume 148 of *ACM International Conference Proceeding Series*, pages 409–416. ACM, 2006. doi: 10.1145/1143844.1143896. URL <https://doi.org/10.1145/1143844.1143896>.
- M. Hutter. Extreme state aggregation beyond mdps. In P. Auer, A. Clark, T. Zeugmann, and S. Zilles, editors, *Algorithmic Learning Theory - 25th International Conference, ALT 2014, Bled, Slovenia, October 8-10, 2014. Proceedings*, volume 8776 of *Lecture Notes in Computer Science*, pages 185–199. Springer, 2014. doi: 10.1007/978-3-319-11662-4\_14. URL [https://doi.org/10.1007/978-3-319-11662-4\\_14](https://doi.org/10.1007/978-3-319-11662-4_14).

- M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=SJ6yPD5xg>.
- P. W. Keller, S. Mannor, and D. Precup. Automatic basis function construction for approximate dynamic programming and reinforcement learning. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, page 449–456, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933832. doi: 10.1145/1143844.1143901. URL <https://doi.org/10.1145/1143844.1143901>.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.
- S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1): 79 – 86, 1951. doi: 10.1214/aoms/1177729694. URL <https://doi.org/10.1214/aoms/1177729694>.
- C. Kwok and D. Fox. Reinforcement learning for sensing strategies. In *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (IEEE Cat. No.04CH37566)*, volume 4, pages 3158–3163 vol.4, 2004. doi: 10.1109/IROS.2004.1389903.
- G. Lambrechts, A. Bolland, and D. Ernst. Informed pomdp: Leveraging additional information in model-based rl, 2023.
- D. Leslie. Reinforcement learning in games. 2004.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In Y. Bengio and Y. LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1509.02971>.
- R. Loynd, R. Fernandez, A. Celikyilmaz, A. Swaminathan, and M. J. Hausknecht. Working memory graphs. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 6404–6414. PMLR, 2020. URL <http://proceedings.mlr.press/v119/loynd20a.html>.
- X. Lu, B. V. Roy, V. Dwaracherla, M. Ibrahimi, I. Osband, and Z. Wen. Reinforcement learning, bit by bit. *CoRR*, abs/2103.04047, 2021. URL <https://arxiv.org/abs/2103.04047>.
- I. Menache, S. Mannor, and N. Shimkin. Basis function adaptation in temporal difference reinforcement learning. *Ann. Oper. Res.*, 134(1):215–238, 2005. doi: 10.1007/s10479-005-5732-z. URL <https://doi.org/10.1007/s10479-005-5732-z>.
- A. W. Moore, J. G. Schneider, and K. Deng. Efficient locally weighted polynomial regression predictions. In *ICML*, 1997.
- A. Mott, D. Zoran, M. Chrzanowski, D. Wierstra, and D. J. Rezende. Towards interpretable reinforcement learning using attention augmented agents. In *NeurIPS*, 2019.
- A. Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997. ISSN 00018678. URL <http://www.jstor.org/stable/1428011>.
- H. Oh and T. Kaneko. Deep recurrent q-network with truncated history. In *2018 Conference on Technologies and Applications of Artificial Intelligence (TAAI)*, pages 34–39, 2018. doi: 10.1109/TAAI.2018.00017.
- OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba. Learning dexterous in-hand manipulation, 2019.
- D. Ormoneit and S. Sen. Kernel-based reinforcement learning. *Mach. Learn.*, 49(2-3):161–178, 2002. doi: 10.1023/A:1017928328829. URL <https://doi.org/10.1023/A:1017928328829>.
- E. Parisotto, H. F. Song, J. W. Rae, R. Pascanu, Ç. Gülçehre, S. M. Jayakumar, M. Jaderberg, R. L. Kaufman, A. Clark, S. Noury, M. Botvinick, N. Heess, and R. Hadsell. Stabilizing transformers for reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 7487–7498. PMLR, 2020. URL <http://proceedings.mlr.press/v119/parisotto20a.html>.
- J. Peng. Efficient memory-based dynamic programming. In A. Prieditis and S. Russell, editors, *Machine Learning Proceedings 1995*, pages 438–446. Morgan Kaufmann, San Francisco (CA), 1995. ISBN 978-1-55860-377-6. doi: <https://doi.org/10.1016/B978-1-55860-377-6.50061-X>. URL <https://www.sciencedirect.com/science/article/pii/B978155860377650061X>.

- M. Petrik. An analysis of laplacian methods for value function approximation in mdps. In M. M. Veloso, editor, *IJCAI 2007, Proceedings of the 20th International Joint Conference on Artificial Intelligence, Hyderabad, India, January 6-12, 2007*, pages 2574–2579, 2007. URL <http://ijcai.org/Proceedings/07/Papers/414.pdf>.
- A. Pritzel, B. Uria, S. Srinivasan, A. P. Badia, O. Vinyals, D. Hassabis, D. Wierstra, and C. Blundell. Neural episodic control. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 2827–2836. PMLR, 2017. URL <http://proceedings.mlr.press/v70/pritzel17a.html>.
- S. Proper and P. Tadepalli. Scaling model-based average-reward reinforcement learning for product delivery. In J. Fürnkranz, T. Scheffer, and M. Spiliopoulou, editors, *Machine Learning: ECML 2006, 17th European Conference on Machine Learning, Berlin, Germany, September 18-22, 2006, Proceedings*, volume 4212 of *Lecture Notes in Computer Science*, pages 735–742. Springer, 2006. doi: 10.1007/11871842\_74. URL [https://doi.org/10.1007/11871842\\_74](https://doi.org/10.1007/11871842_74).
- A. Raichuk, P. Stanczyk, M. Orsini, S. Girgin, R. Marinier, L. Hussenot, M. Geist, O. Pietquin, M. Michalski, and S. Gelly. What matters for on-policy deep actor-critic methods? a large-scale study. In *ICLR*, 2021.
- S. Ritter, R. Faulkner, L. Sartran, A. Santoro, M. Botvinick, and D. Raposo. Rapid task-solving in novel environments. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL [https://openreview.net/forum?id=F-mvpFpn\\_0q](https://openreview.net/forum?id=F-mvpFpn_0q).
- S. S. Ruan, G. Comanici, P. Panangaden, and D. Precup. Representation discovery for mdps using bisimulation metrics. In *AAAI*, 2015.
- A. Schaefer, S. Udluft, and H.-G. Zimmermann. A recurrent control neural network for data efficient reinforcement learning. *2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning*, pages 151–157, 2007.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- D. Sejdinovic, B. Sriperumbudur, A. Gretton, and K. Fukumizu. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, 41(5):2263–2291, 2013. ISSN 00905364, 21688966. URL <http://www.jstor.org/stable/23566550>.
- S. P. Singh, D. J. Litman, M. J. Kearns, and M. A. Walker. Optimizing dialogue management with reinforcement learning: Experiments with the njfun system. *J. Artif. Intell. Res.*, 16:105–133, 2002. doi: 10.1613/jair.859. URL <https://doi.org/10.1613/jair.859>.
- I. Sorokin, A. Seleznev, M. Pavlov, A. Fedorov, and A. Ignateva. Deep attention recurrent q-network. *ArXiv*, abs/1512.01693, 2015.
- B. K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. R. G. Lanckriet. On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6(none):1550 – 1599, 2012. doi: 10.1214/12-EJS722. URL <https://doi.org/10.1214/12-EJS722>.
- J. Subramanian and A. Mahajan. Approximate information state for partially observed systems. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 1629–1636, 2019. doi: 10.1109/CDC40024.2019.9029898.
- J. Subramanian, A. Sinha, R. Seraj, and A. Mahajan. Approximate information state for approximate planning and reinforcement learning in partially observed systems. *CoRR*, abs/2010.08843, 2020. URL <https://arxiv.org/abs/2010.08843>.
- R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, USA, 1998. ISBN 0-262-19398-1. URL <http://www.cs.ualberta.ca/%7Esutton/book/ebook/the-book.html>.
- R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pages 1057–1063. The MIT Press, 1999. URL <https://proceedings.neurips.cc/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf>.
- Y. Tang, D. Nguyen, and D. Ha. Neuroevolution of self-interpretable agents. In C. A. C. Coello, editor, *GECCO '20: Genetic and Evolutionary Computation Conference, Cancún Mexico, July 8-12, 2020*,

pages 414–424. ACM, 2020. doi: 10.1145/3377930.3389847. URL <https://doi.org/10.1145/3377930.3389847>.

E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.

J. N. Tsitsiklis and B. V. Roy. Feature-based methods for large scale dynamic programming. *Mach. Learn.*, 22(1-3):59–94, 1996. doi: 10.1023/A:1018008221616. URL <https://doi.org/10.1023/A:1018008221616>.

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.

R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 2004.

K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.

X. Xu, T. Xie, D. Hu, and X. Lu. Kernel least-squares temporal difference learning. In *Computational intelligence and neuroscience*, 2006.

V. F. Zambaldi, D. Raposo, A. Santoro, V. Bapst, Y. Li, I. Babuschkin, K. Tuyls, D. P. Reichert, T. P. Lillicrap, E. Lockhart, M. Shanahan, V. Langston, R. Pascanu, M. M. Botvinick, O. Vinyals, and P. W. Battaglia. Deep reinforcement learning with relational inductive biases. In *ICLR*, 2019.

## Checklist

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes]
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Not Applicable]
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Not Applicable]
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. [Yes]
  - (b) Complete proofs of all theoretical results. [Yes]
  - (c) Clear explanations of any assumptions. [Yes]
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Not Applicable]
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Not Applicable]
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes]
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Not Applicable]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (a) Citations of the creator If your work uses existing assets. [Yes]
  - (b) The license information of the assets, if applicable. [Not Applicable]
  - (c) New assets either in the supplemental material or as a URL, if applicable. Not Applicable]
  - (d) Information about consent from data providers/curators. [Not Applicable]
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
  - (a) The full text of instructions given to participants and screenshots. [Not Applicable]
  - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Not Applicable]
  - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Not Applicable]

## Appendix

### A Proof for Theorem 3.5

For readability we will restate the theorem statement

**Theorem A.1.** *For any time  $t$ , any realisation  $s_t$  of  $S_t$ ,  $a_t$  of  $A_t$ , let  $h_t = (s_{1:t}, a_{1:t-1})$ , and  $z_t = \sigma_t(h_t)$ .*

*The worst case difference between  $V^*$  and  $V_t^\pi$  is bounded as:*

$$\Delta \leq 2 \frac{\varepsilon + \gamma \delta \kappa_{\mathfrak{F}}(\hat{V}, \hat{f})}{1 - \gamma}, \quad (29)$$

where,  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f}) = \sup_{z,a} \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z, a)))$ . and  $\rho_{\mathfrak{F}}(\cdot)$  is the Minkowski functional associated with the IPM  $d_{\mathfrak{F}}$  as defined in (3).

*Proof.* For this proof we will use the following convention: For a generic history  $h_t \in \mathcal{H}_t$ , we assume that  $h_t = (s_{1:t}, a_{1:t-1})$ , moreover, note that  $z_t = \sigma_t(h_t)$ .

Now from (3.1), and Definition 3.3 for any  $a_t, s_t, z_t$ :

$$\begin{aligned} & \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| r(s_t, a_t) - \hat{r}(z_t, a_t) \right| \leq \varepsilon. \\ & \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| \sum_{s_{t+1} \in \mathcal{S}} \left( P(s_{t+1}|s_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) - \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right) \right| \leq \delta \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z_t, a_t))). \end{aligned} \quad (30)$$

Now using triangle inequality we get:

$$\|V^*(s_t) - V_t^\pi(h_t)\|_\infty \leq \underbrace{\|V^*(s_t) - \hat{V}(z_t)\|_\infty}_{\text{term 1}} + \underbrace{\|V_t^\pi(h_t) - \hat{V}(z_t)\|_\infty}_{\text{term 2}}, \quad (31)$$

We will now proceed by bounding terms 1 and 2 separately

#### Bounding term 1:

$$\|V^*(s_t) - \hat{V}(z_t)\|_\infty \leq \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ Q^*(s_t, a_t) - \hat{Q}(z_t, a_t) \right] \right|, \quad (32)$$

For any action  $a_t$  we have:

$$\begin{aligned} & \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ Q^*(s_t, a_t) - \hat{Q}(z_t, a_t) \right] \right| = \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ r(s_t, a_t) + \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V^*(s_{t+1}) \right. \right. \\ & \quad \left. \left. - \hat{r}(z_t, a_t) - \gamma \sum_{s_{t+1} \in \mathcal{S}} \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right] \right| \\ & \stackrel{(a)}{\leq} \varepsilon + \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V^*(s_{t+1}) - \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right| \\ & \quad + \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) - \gamma \sum_{s_{t+1} \in \mathcal{S}} \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right| \\ & \stackrel{(b)}{\leq} \varepsilon + \gamma \|V^*(s_t) - \hat{V}(z_t)\|_\infty + \gamma \delta \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z_t, a_t))), \end{aligned}$$

where (a) from triangle inequality and (b) is due to (30). Now defining  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f}) = \sup_{z,a} \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z, a)))$ , and substituting the above result in (32) we get

$$\|V^*(s_t) - \hat{V}(z_t)\|_\infty \leq \frac{\varepsilon + \gamma \delta \kappa_{\mathfrak{F}}(\hat{V}, \hat{f})}{1 - \gamma}. \quad (33)$$



**Bounding term 2:**

$$\|V_t^\pi(h_t) - \hat{V}(z_t)\|_\infty \leq \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ Q_t^\pi(h_t, a_t) - \hat{Q}(z_t, a_t) \right] \right|, \quad (34)$$

For any action  $a_t$ , we have:

$$\begin{aligned} \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ Q^\pi(h_t, a_t) - \hat{Q}(z_t, a_t) \right] \right| &= \max_{h \in \mathcal{H}_t} \left| \max_{a_t \in \mathcal{A}} \left[ r(s_t, a_t) + \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V_{t+1}^\pi(h_{t+1}) - (\hat{r}(z_t, a_t) \right. \right. \\ &\quad \left. \left. + \gamma \sum_{s_{t+1} \in \mathcal{S}} \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right] \right|, \\ &\stackrel{(a)}{\leq} \epsilon + \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) V_{t+1}^\pi(h_{t+1}) - \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right| \\ &\quad + \max_{h \in \mathcal{H}_t, a_t \in \mathcal{A}} \left| \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) - \gamma \sum_{s_{t+1} \in \mathcal{S}} \hat{P}(s_{t+1}|z_t, a_t) \hat{V}(\hat{f}(s_{t+1}, z_t, a_t)) \right|, \\ &\stackrel{(b)}{\leq} \epsilon + \gamma \| (V^\pi(h_t) - \hat{V}(z_t)) \|_\infty + \gamma \delta \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z_t, a_t))), \end{aligned}$$

where (a) is from triangle inequality, (b) is due to (30), with  $\kappa_{\mathfrak{F}}(\hat{V}, \hat{f}) = \sup_{z, a} \rho_{\mathfrak{F}}(\hat{V}(\hat{f}(\cdot, z, a)))$ , and substituting the above result in (34) we get

$$\|V_t^\pi(h_t) - \hat{V}(z_t)\|_\infty \leq \frac{\epsilon + \gamma \delta \kappa_{\mathfrak{F}}(\hat{V}, \hat{f})}{1 - \gamma}. \quad (35)$$

The final result follows by adding (33) and (35).  $\square$

## B Proof for Corollary 3.6

**Lemma B.1.** *If  $\hat{V}$  is the optimal value function of the MDP  $\hat{\mathcal{M}}$  induced by the process  $\{Z_t\}_{t \geq 0}$ , then*

$$\text{span}(\hat{V}) \leq \frac{\text{span}(\hat{r})}{1 - \gamma}. \quad (36)$$

*Proof.* The result follows by observing that the per-step reward  $\hat{r}(Z_t, A_t) \in [\min(\hat{r}), \max(\hat{r})]$ . Therefore  $\max(\hat{V}) \leq \max(\hat{r})$  and  $\min(\hat{V}) \geq \min(\hat{r})$ .  $\square$

**Corollary B.2.** *If the function class  $\mathfrak{F}$  is  $\mathfrak{F}^{TV}$ , then  $\Delta$  defined in (14) is upper bounded as:*

$$\Delta \leq \frac{2\epsilon}{1 - \gamma} + \frac{\gamma \delta \text{span}(\hat{r})}{(1 - \gamma)^2}, \quad (37)$$

*Proof.* From Section 3.1 we know that for the Total variation distance  $\rho_{\mathfrak{F}^{TV}}(\hat{V}) = \text{span}(\hat{V})$  and  $\kappa(\hat{f}) = 1$ . The result in the corollary then follows from Lemma B.1.  $\square$

## C Proof for Corollary 3.7

**Definition C.1.** For any Lipschitz function  $f : (\mathcal{Z}, d_Z) \rightarrow (\mathcal{R}, |\cdot|)$ , and probability measures  $\nu_1$ , and  $\nu_2$  on  $(\mathcal{Z}, d_Z)$

$$\left| \int_{\mathcal{Z}} f d\nu_1 - \int_{\mathcal{Z}} f d\nu_2 \right| \leq \|f\|_L \cdot d_{\mathfrak{F}^w}(\nu_1, \nu_2) \leq L_f d_{\mathfrak{F}^w}(\nu_1, \nu_2), \quad (38)$$

where  $L_f$  is the Lipschitz constant of  $f$  and  $d_{\mathfrak{F}^w}$  is the Wasserstein distance.

**Definition C.2.** Let  $d$  be a metric on the AIS/Feature space  $\mathcal{Z}$ . The MDP  $\hat{\mathcal{M}}$  induced by the process  $\{Z_t\}_{t \geq 0}$  is said to be  $(L_{\hat{r}}, L_{\hat{P}})$  - Lipschitz if for any  $Z_1, Z_2 \in \mathcal{Z}$ , the reward  $\hat{r}$  and transition  $\hat{P}$  of  $\hat{\mathcal{M}}$  satisfy the following:

$$\left| r(Z_1, A) - r(Z_2, A) \right| \leq L_{\hat{r}} d(Z_1, Z_2), \quad (39)$$

$$d_{\mathfrak{S}^w}(\hat{P}(\cdot|Z_1, A), \hat{P}(\cdot|Z_2, A)) \leq L_{\hat{P}} d(Z_1, Z_2), \quad (40)$$

where  $d_{\mathfrak{S}^w}$  is the Wasserstein or the Kantorovitch-Rubinstein distance.

**Lemma C.3.** Let  $\hat{V} : \mathcal{Z} \rightarrow \mathcal{R}$  be  $L_{\hat{V}}$  continuous. Define:

$$\hat{Q}(z, a) = \hat{r}(z, a) + \gamma \sum_{s'} \hat{P}(s'|z, a) \hat{V}(f(s', z, a)).$$

Then  $\hat{Q}$  is  $(L_{\hat{r}} + \gamma L_{\hat{V}} L_{\hat{f}} L_{\hat{P}})$ -Lipschitz continuous.

*Proof.* For any action  $a$

$$\left| \hat{Q}(z_1, a) - \hat{Q}(z_2, a) \right| \stackrel{(a)}{\leq} \left| \hat{r}(z_1, a) - \hat{r}(z_2, a) \right| + \gamma \left| \sum_{s'} \hat{P}(s'|z_1, a) \hat{V}(f(s', z_1, a)) - \sum_{s'} \hat{P}(s'|z_2, a) \hat{V}(f(s', z_2, a)) \right|, \quad (41)$$

$$\stackrel{(b)}{\leq} (L_{\hat{r}} + \gamma L_{\hat{V}} L_{\hat{f}} L_{\hat{P}}) d(z_1, z_2), \quad (42)$$

where (a) due to triangle inequality, and (b) follows from Definition C.1, Definition C.2, and because  $\|a \circ b\|_L \leq \|a\|_L \cdot \|b\|_L$ .  $\square$

**Lemma C.4.** Let  $\hat{Q} : \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{R}$  be  $L_{\hat{Q}}$ - Lipschitz continuous, Define

$$\hat{V}(z) = \max_{a_t \in \mathcal{A}} \hat{Q}(z, a).$$

Then  $\hat{V}$  is  $L_{\hat{Q}}$  Lipschitz

*Proof.* Consider  $z_1, z_2 \in \mathcal{Z}$ , and let  $a_1$  and  $a_2$  denote the corresponding optimal action. Then,

$$\hat{V}(z_1) - \hat{V}(z_2) = \hat{Q}(z_1, a_1) - \hat{Q}(z_2, a_2), \quad (43)$$

$$\stackrel{(a)}{\leq} \hat{Q}(z_1, a_2) - \hat{Q}(z_2, a_2), \quad (44)$$

$$\stackrel{(b)}{\leq} L_{\hat{Q}} d(z_1, z_2), \quad (45)$$

By symmetry,

$$\hat{V}(z_2) - \hat{V}(z_1) \leq L_{\hat{Q}} d(z_1, z_2).$$

Therefore,

$$\left| \hat{V}(z_1) - \hat{V}(z_2) \right| \leq L_{\hat{Q}} d(z_1, z_2). \quad \square$$

**Lemma C.5.** Consider the following dynamic program defined in (9):<sup>2</sup>

$$\hat{Q}_t(z_t, a_t) = \hat{r}(z_t, a_t) + \gamma \sum_{s_t \in \mathcal{S}} \hat{P}(s_t|z_t, a_t) \hat{V}(f(z_t, s_t, a_t)), \quad \forall z \in \mathcal{Z}, a \in \mathcal{A}$$

$$\hat{V}_t(z_t) = \max_{a \in \mathcal{A}} \hat{Q}_t(z_t, a), \quad \forall z \in \mathcal{Z}$$

Then at any time  $t$ , we have:

$$L_{\hat{V}_{t+1}} = L_{\hat{r}} + \gamma L_{\hat{P}} L_{\hat{f}} L_{\hat{V}_t}.$$

<sup>2</sup>We have added  $t$  as a subscript to denote the computation time *i.e.*, the time at which the respective function is updated.

*Proof.* We prove this by induction. At time  $t = 1$   $\hat{Q}_1(z, a) = \hat{r}(z, a)$ , therefore  $L_{\hat{Q}_1} = L_{\hat{r}}$ . Then according to Lemma C.4,  $\hat{V}_1$  is Lipschitz with Lipschitz constant  $L_{\hat{V}_1} = L_{\hat{Q}_1} = L_{\hat{r}}$ . This forms the basis of induction. Now assume that at time  $t$ ,  $\hat{V}_t$  is  $L_{\hat{V}_t}$ -Lipschitz. By Lemma C.3  $\hat{Q}_{t+1}$  is  $L_{\hat{r}} + \gamma L_{\hat{f}}, L_{\hat{P}} L_{\hat{V}_t}$ . Therefore by Lemma C.4,  $\hat{V}_{(t+1)}$  is Lipschitz with constant:

$$L_{\hat{V}_{t+1}} = L_{\hat{r}} + \gamma L_{\hat{f}} L_{\hat{P}} L_{\hat{V}_t}.$$

□

**Theorem C.6.** *Given any  $(L_{\hat{r}}, L_{\hat{P}})$ -Lipschitz MDP, if  $\gamma L_{\hat{P}} L_{\hat{f}} \leq 1$ , then the infinite horizon  $\gamma$ -discounted value function  $\hat{V}$  is Lipschitz continuous with Lipschitz constant*

$$L_{\hat{V}} = \frac{L_{\hat{r}}}{1 - \gamma L_{\hat{f}} L_{\hat{P}}}.$$

*Proof.* Consider the sequence of  $L_t = L_{\hat{V}_t}$  values. For simplicity write  $\alpha = \gamma L_{\hat{P}} L_{\hat{f}}$ . Then the sequence  $\{L_t\}_{t \geq 1}$  is given by :  $L_1 = L_{\hat{r}}$  and for  $t \geq 1$ ,

$$L_{t+1} = L_{\hat{r}} + \alpha L_t,$$

Therefore,

$$L_t = L_{\hat{r}} + \alpha L_{\hat{r}} + \dots + \alpha_{t+1} = \frac{1 - \alpha^t}{1 - \alpha} L_{\hat{r}}.$$

This sequence converges if  $|\alpha| \leq 1$ . Since  $\alpha$  is non-negative, this is equivalent to  $\alpha \leq 1$ , which is true by hypothesis. Hence  $L_t$  is a convergent sequence. At convergence, the limit  $L_{\hat{V}}$  must satisfy the fixed point of the recursion relationship introduced in Lemma C.5, hence,

$$L_{\hat{V}} = L_{\hat{r}} + \gamma L_{\hat{f}} L_{\hat{P}} L_{\hat{V}}.$$

Consequently, the limit is equal to,

$$L_{\hat{V}} = \frac{L_{\hat{r}}}{1 - \gamma L_{\hat{f}} L_{\hat{P}}}.$$

□

**Corollary C.7.** *If  $\gamma L_{\hat{P}} L_{\hat{f}} \leq 1$  and the function class  $\mathfrak{F}$  is  $\mathfrak{F}^W$ , then  $\Delta$  as defined in (14) is upper bounded as:*

$$\Delta \leq \frac{2\epsilon}{(1 - \gamma)} + \frac{2\gamma\delta L_{\hat{r}}}{(1 - \gamma)(1 - \gamma L_{\hat{f}} L_{\hat{P}})}, \quad (46)$$

*Proof.* The proof follows from the observation that for  $d_{\mathfrak{F}^W}, \rho_{\mathfrak{F}^W} = L_{\hat{V}}$ , and then using the result from Theorem C.6. □

## D Algorithmic Details

### D.1 Choice of an IPM:

#### D.1.1 MMD

One advantage of choosing  $d_{\mathfrak{F}}$  as the MMD distance is that unlike the Wasserstein distance, its computation does not require solving an optimisation problem. Another advantage is that we can leverage some of their properties to further simplify our computation, as follows:

**Proposition D.1** (Theorem 22 [Sejdinovic et al., 2013]). *Let  $\mathcal{X} \subseteq \mathcal{R}^m$ , and  $d_{\mathcal{X},p} : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{R}$  be a metric given by  $d_{\mathcal{X},p}(x, x') = \|x - x'\|_2^p$ , for  $p \in (0, 2]$ . Let  $k_p : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{R}$  be any kernel given:*

$$k_p(x, x') = \frac{1}{2}(d_{\mathcal{X},p}(x, x_0) + d_{\mathcal{X},p}(x', x_0) - d_{\mathcal{X},p}(x, x')), \quad (47)$$

where  $x_0 \in \mathcal{X}$  is arbitrary, and let  $\mathcal{U}_p$  be a RKHS kernel with kernel  $k_p$  and  $\mathfrak{F}_p = \{f \in \mathcal{U}_p : \|f\|_{\mathcal{U}_p} \geq 1\}$ . Then for any distributions  $\nu_1, \nu_2 \in \Delta\mathcal{X}$ , the IPM can be expressed as:

$$d_{\mathfrak{F}}(\nu_1, \nu_2) = \left( \mathbb{E}[d_{\mathcal{X},p}(X_1, W_1)] - \frac{1}{2}\mathbb{E}[d_{\mathcal{X},p}(X_1, X_2)] - \frac{1}{2}\mathbb{E}[d_{\mathcal{X},p}(W_1, W_2)] \right)^{\frac{1}{2}}, \quad (48)$$

where  $X_1, X_2$ , and  $W_1, W_2$  are i.i.d. samples from  $\nu_1$  and  $\nu_2$  respectively.

The main implication of Proposition D.1 is that, instead of using (48), for  $p \in (0, 2]$  we can use the following as a surrogate for  $d_{\mathfrak{F}_p}$ :

$$\int_{\mathcal{X}} \int_{\mathcal{X}} \|x_1 - w_1\|_2^p \nu_1(dx_1) \nu_2(dw_1) - \frac{1}{2} \int_{\mathcal{X}} \int_{\mathcal{X}} \|w_1 - w_2\|_2^p \nu_2(dw_1) \nu_2(dw_2). \quad (49)$$

Moreover, according to Sriperumbudur et al. [2012] for  $n$  identically and independently distributed (i.i.d) samples  $\{X_i\}_{i=0}^n \sim \nu_1$  an unbiased estimator of (49) is given as:

$$\frac{1}{n} \sum_{i=1}^n \int_{\mathcal{X}} \|X_i - w_1\|_2^p \nu_1(dw_1) - \frac{1}{2} \int_{\mathcal{X}} \int_{\mathcal{X}} \|w_1 - w_2\|_2^p \nu_1(dw_1) \nu_2(dw_2). \quad (50)$$

We implement a simplified version of the surrogate loss in (50) as follows:

**Proposition D.2** ([Subramanian et al., 2020]). *Given the setup in Proposition D.1 and  $p = 2$ , Let  $\nu_2(\zeta)$  be a parametric distribution with mean  $m$  and let  $X \sim \nu_1$ , then the gradient  $\nabla_{\zeta}(m_{\zeta} - 2X)^{\top} m_{\zeta}$  is an unbiased estimator of  $\nabla_{\zeta} d_{\mathfrak{F}_2}(\alpha, \nu_{\zeta})^2$*

*Proof.* Let  $X_1, X_2 \sim \nu_1$ , and  $W_1, W_2 \sim \nu_2(\zeta)$

$$\therefore \nabla_{\zeta} d_{\mathfrak{F}_2}(\nu_1, \nu_2(\zeta))^2 = \nabla_{\zeta} \left[ \mathbb{E}\|X_1 - W_1\|_2^2 - \frac{1}{2}\mathbb{E}\|X_1 - X_2\|_2^2 - \frac{1}{2}\mathbb{E}\|W_1 - W_2\|_2^2 \right], \quad (51)$$

$$\stackrel{(a)}{=} \nabla_{\zeta} \left[ \mathbb{E}\|W_1\|_2^2 - 2\mathbb{E}\|X_1\|^{\top} \mathbb{E}\|W_1\| \right], \quad (52)$$

where (a) follows from the fact that  $X$  does not depend on  $\zeta$ , which simplifies the implementation of the MMD distance.  $\square$

In this way we can simplify the computation of  $d_{\mathfrak{F}}$  using a parametric stochastic kernel approximator and MMD metric.

Note that when are trying to approximate a continuous distribution we can readily use the loss function (52) as long as the mean  $m_{\zeta}$  of  $\nu_2(\zeta)$  is given in closed form. The AIS loss is then given as:

$$\mathcal{L}_{\text{AIS}}(\zeta) = \frac{1}{T} \sum_{t=0}^T \left( \lambda(f_{\hat{r}}(Z_t, A_t; \zeta) - r(S_t, A_t))^2 + (1 - \lambda)(m_{\zeta}^{S_t} - 2S_t)^{\top} m_{\zeta}^{S_t} \right), \quad (53)$$

where  $m_{\zeta}^{S_t}$  is obtained using the from the transition approximator, i.e., the mapping  $f_{\hat{P}}(\zeta) : \mathcal{Z} \times \mathcal{A} \rightarrow \mathcal{R}$ .

### D.1.2 Wasserstein Distance

The the KL-divergence between two densities  $\nu_1$  and  $\nu_2$  on for any  $X \in \mathcal{X} \subset \mathcal{R}^m$  is defined as:

$$d_{\text{KL}}(\nu_1 \parallel \nu_2) = \int_{\mathcal{X}} \log(\nu_1(x)) \nu_1(dx) - \int_{\mathcal{X}} \log(\nu_2(x)) \nu_1(dx) \quad (54)$$

Moreover, if  $\mathcal{X}$  is bounded space with diameter  $D$ , then the relation between the Wasserstein distance  $d_{\mathfrak{F}^w}$ , Total variation distance  $d_{\mathfrak{F}^{\text{TV}}}$ , and the KL divergence is given as :

$$d_{\mathfrak{F}^w}(\nu_1, \nu_2) \leq D d_{\mathfrak{F}^{\text{TV}}}(\nu_1, \nu_2) \stackrel{(a)}{\leq} \sqrt{2d_{\text{KL}}(\nu_1 \parallel \nu_2)}, \quad (55)$$

where, (a) follows from the Pinsker’s inequality. Note that in (19) we use  $d_{\frac{2}{3}}^2$ . Therefore, we can use a (simplified) KL-divergence based surrogate objective given as:

$$\int_{\mathcal{X}} \log(\nu_2(x; \zeta)) \nu_1(dx), \quad (56)$$

where we have dropped the terms which do not depend on  $\zeta$ . Note that the above expression is same as the cross entropy between  $\nu_1$  and  $\nu_2$  which can be effectively computed using samples. In particular, if we get  $T$  i.i.d samples from  $\nu_1$ , then,

$$\frac{1}{T} \sum_{i=0}^T \log(\nu_2(x_i; \zeta)) \quad (57)$$

is an unbiased estimator of  $\int_{\mathcal{X}} \log(\nu_2(x; \zeta)) \nu_1(dx)$ .

The KL divergence based AIS loss is then given as:

$$\mathcal{L}_{\text{AIS}}(\zeta) = \frac{1}{T} \sum_{t=0}^T \left( \lambda (f_{\hat{r}}(Z_t, A_t; \zeta) - r(S_t, A_t))^2 + (1 - \lambda) \log(\hat{P}(S_t; \zeta)) \right), \quad (58)$$

## E Convergence Analysis

In this section we will discuss the convergence of the AIS-based policy gradient in Section 4.2 as well as Actor-Critic algorithm presented in the previous subsection. The proof of convergence relies on multi-timescale stochastic approximation Borkar [2008] under conditions similar to the standard conditions for convergence of policy gradient algorithms with function approximation stated below, therefore it would suffice to provide a proof sketch.

- Assumption E.1.** 1. The values of step-size parameters  $b, d$  and  $c$  (for the actor critic algorithm) are set such that the timescales of the updates for  $\zeta, \xi$ , and  $\vartheta$  (for Actor-Critic algorithm) are separated, *i.e.*,  $b_t \gg d_t$ , and for the Actor-Critic algorithm  $b_t \gg c_t \gg d_t$ ,  $\sum_i b_i = \infty$ ,  $\sum_i b_i^2 < \infty$ ,  $\sum_i c_i = \infty$ ,  $\sum_i c_i^2 < \infty$ ,  $\sum_i d_i = \infty$  and  $\sum_i d_i^2 < \infty$ ,  $\lim_{i \rightarrow \infty} \frac{d_i}{b_i} = 0$  and  $\lim_{i \rightarrow \infty} \frac{c_i}{d_i} = 0$ ,
2. The parameters  $\zeta, \xi$  and  $\vartheta$  (for Actor-Critic algorithm) lie in a convex, compact and closed subset of Euclidean spaces.
  3. The gradient  $\nabla_{\zeta} \mathcal{L}_{\text{AIS}}$  is Lipschitz in  $\zeta_t$ , and  $\hat{\nabla}_{\xi} J(\xi, \zeta)$  is Lipschitz in  $\xi_t$ , and  $\zeta_t$ . Whereas for the Actor-Critic algorithm the gradient of the TD loss  $\nabla_{\vartheta} \mathcal{L}_{\text{TD}}(\zeta, \xi, \vartheta)$  and the policy gradient  $\hat{\nabla}_{\xi} J(\zeta, \xi, \vartheta)$  is Lipschitz in  $(\zeta_t, \xi_t, \vartheta_t)$ .
  4. All the estimates of all the gradients  $\nabla_{\zeta} \mathcal{L}_{\text{AIS}}$ ,  $\nabla_{\xi} J(\xi, \zeta)$ ,  $\nabla_{\vartheta} \mathcal{L}_{\text{TD}}(\zeta, \xi, \vartheta)$  and are unbiased with bounded variance<sup>3</sup>.

- Assumption E.2.** 1. The ordinary differential equation (ODE) corresponding to (23) is locally asymptotically stable.
2. The ODEs corresponding to (22) is globally asymptotically stable.
  3. For the Actor-Critic algorithm, the ODE corresponding to (28b) is globally asymptotically stable and has a fixed point which is Lipschitz in  $\xi$ .

**Theorem E.3.** *Under assumption E.1 and E.2, along any sample path, almost surely we have the following:*

1. *The iteration for  $\zeta$  in (22) converges to an AIS generator that minimises the  $\mathcal{L}_{\text{AIS}}$ .*
2. *The iteration for  $\xi$  in (23) converges to a local maximum of the performance  $J(\zeta^*, \xi)$  where  $\zeta^*$ , and  $\vartheta^*$  (for Actor-Critic) are the converged value of  $\zeta, \vartheta$ .*
3. *For the Actor-Critic algorithm the iteration for  $\vartheta$  in (28b) converges to critic that minimises the error with respect to the true value function.*

---

<sup>3</sup>This assumption is only satisfied in tabular MDPs.

*Proof.* The proof for this theorem follows the technique used in [Leslie, 2004, Borkar, 2008]. Due to the specific choice of learning rate the AIS-generator is updated at a faster time-scale than the actor, therefore it is “quasi static” with respect to the actor while the actor observes a “nearly equilibrated” AIS generator. Similarly in the case of the Actor-Critic algorithm the AIS generator observes a stationary critic and actor, whereas the critic and actor see “nearly equilibrated” AIS generator. The Martingale difference condition (A3) of Borkar [2008] is satisfied due to Item 4 in assumption E.1. As such since our algorithm satisfies all the four conditions by [Leslie, 2004, page35], [Borkar, 1997, Theorem 23], the result then follows by combining the theorem on [Leslie, 2004, page 35][Borkar, 2008, Theorem 23] and [Borkar, 1997, Theorem 2.2].  $\square$

## F Experimental Details

Common	Optimiser	Adam
	Discount Factor $\gamma$	0.99
	Initial standard deviation for the policy	0.0
	PPO-Epochs	12
	Clipping Coefficient	0.2
	Entropy-Regulariser	0
	Batch Size	512
AIS generator	Episode Length	2048
	History Compressor	GRU
	Hidden layer dimension	256
	Step size	1.5e-3
Actor	$\lambda$	0.3
	Step size	3.5e-4
	No of hidden layers	1
	Hidden layer Dimension	32

Table 1: Hyperparameters

### F.1 Environments

Our algorithms are evaluated on MuJoCo [Todorov et al., 2012, mujoco-py version 2.0.2.9 ] via OpenAI gym [Brockman et al., 2016, version 0.17.1] interface, using the v2 environments. The environment, state-space, action space, and reward function are not modified or pre-processed in any way for easy reproducibility and fair comparison with previous results. Each environment runs for a maximum of 2048 time steps or until some termination condition and has a multi-dimensional action space with values in the range of (-1, 1), except for Humanoid which uses the range of (-0.4, 0.4).

### F.2 Hyper-parameters

Table 1 contains all the hyper-parameters used in our experiments. Both the policy and AIS networks are trained with Adam optimiser [Kingma and Ba, 2015], with a batch size of 512. We follow Raichuk et al. [2021]’s recommended protocol for training on-policy policy based methods, and perform 12 PPO updates after every policy evaluation subroutine. To ensure separation of time-scales the step size of the AIS generator and the policy network is set to  $1.5e^{-3}$  and  $3.5e^{-4}$  respectively. Hyper-parameters of our approach are searched over a grid of values, but an exhaustive grid search is not carried out due to prohibitive computational cost. We start with the recommended hyper-parameters for the baseline implementations and tune them further around promising values by an iterative process of performing experiments and observing results.

For the state-based RNN baseline we have tuned the learning rate over a grid of values starting from  $1e^{-4}$  to  $4e^{-4}$  and settled on  $3.5e^{-4}$  as it achieved the best performance. Similarly the hidden layer size set to 256 as it is observed to achieve best performance. For the feed-forward baselines we use the implementation by OpenAI baselines [Dhariwal et al., 2017] with their default hyper-parameters.

### F.3 Modifications to baselines and their hyper-parameters

Note that the methods in [Hafner et al., 2020, 2019] are designed for pixel-based control tasks and cannot be readily used for continuous control tasks in this paper. To help them process real-valued state vectors, we replace the convolutional and deconvolutional layers in their architectures by fully connected layers. We observed that feed-forward layers of size 256 for PlaNET [Hafner et al., 2019], and 300 for Dreamer [Hafner et al., 2020] produced the best results for both the methods. For PlaNET [Hafner et al., 2019] we used the default hyper-parameters and varied the learning rate over a grid of values starting from 1e-4 to 4e-4. We observed that this method achieved best performance at the learning rate of 3e-3. For Dreamer [Hafner et al., 2020] we used the default hyper-parameters and varied the learning rate over a grid of values starting from 6e-4, 6e-5 and 6e-5 respectively to 1e-3, 1e-4, 1e-4 respectively. We observed that this method achieved the best performance at 7.5e-4, 9e-5 and 8.5e-5 respectively.

### F.4 Type of MMDs

The MMD distance given by (52) in Appendix D.1.1, can be computed using different types of characteristic kernels (for a detailed review see [Sriperumbudur et al., 2012, Fukumizu et al., 2009, Sejdinovic et al., 2013]). In this paper we consider computing (52) using the Laplace, Gaussian and energy distance kernels. The performance of the proposed approach under different kernels is shown in Figure 4. It can be observed that for the continuous control tasks in the MuJoCo suite, the energy distance yields better performance, and therefore we implement Equation (52) using the energy distance for all the experiments.

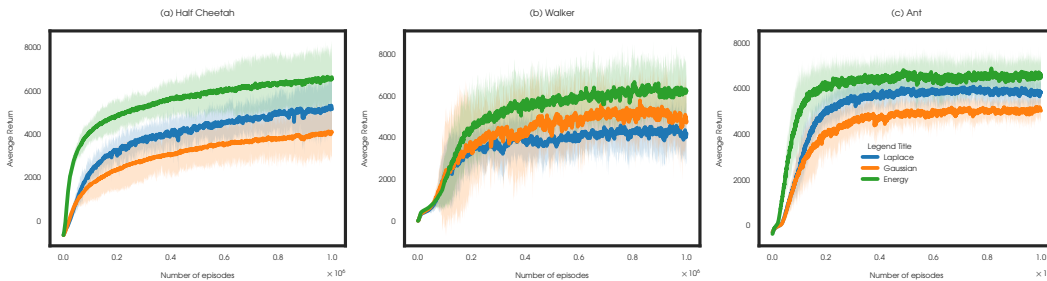


Figure 4: Comparison of different MMDs, averaged over 50 runs with shaded regions representing the quantiles

### F.5 MMD vs KL

Next, we compare the performance of our method under MMD (energy distance)-based AIS loss in (20) and KL-based AIS loss given in (21). From Figure 5, one can observe that for the Mujoco tasks, MMD-based loss leads to better performance.

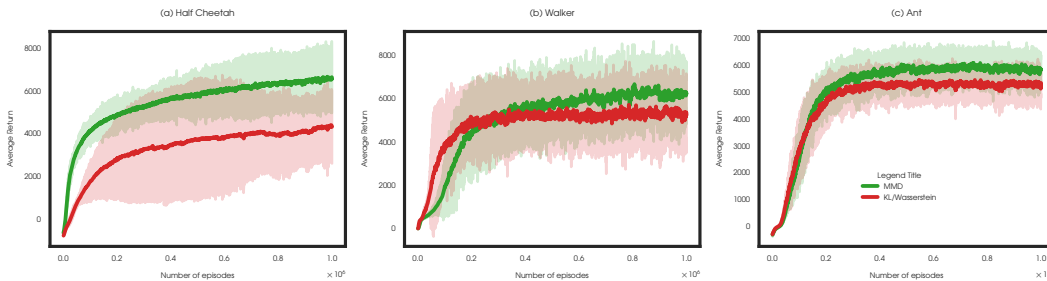


Figure 5: Comparison of Wasserstein vs MMDs, averaged over 50 runs with shaded regions representing the quantiles