

---

# Sample Complexity of Reinforcement Learning using Linearly Combined Model Ensembles

---

Aditya Modi<sup>1</sup>    Nan Jiang<sup>2</sup>

<sup>1</sup>University of Michigan Ann Arbor

Ambuj Tewari<sup>1</sup>    Satinder Singh<sup>1</sup>

<sup>2</sup>University of Illinois at Urbana-Champaign

## Abstract

Reinforcement learning (RL) methods have been shown to be capable of learning intelligent behavior in rich domains. However, this has largely been done in simulated domains without adequate focus on the process of building the simulator. In this paper, we consider a setting where we have access to an ensemble of pre-trained and possibly inaccurate simulators (models). We approximate the real environment using a state-dependent linear combination of the ensemble, where the coefficients are determined by the given state features and some unknown parameters. Our proposed algorithm provably learns a near-optimal policy with a sample complexity polynomial in the number of unknown parameters, and incurs no dependence on the size of the state (or action) space. As an extension, we also consider the more challenging problem of model selection, where the state features are unknown and can be chosen from a large candidate set. We provide exponential lower bounds that illustrate the fundamental hardness of this problem, and develop a provably efficient algorithm under additional natural assumptions.

## 1 INTRODUCTION

Reinforcement learning methods with deep neural networks as function approximators have recently demonstrated prominent success in solving complex and observations rich tasks like games (Mnih et al., 2015; Silver et al., 2016), simulated control problems (Todorov et al., 2012; Lillicrap et al., 2015; Mordatch et al., 2016) and

a range of robotics tasks (Christiano et al., 2016; Tobin et al., 2017). A common aspect in most of these success stories is the use of simulation. Arguably, given a simulator of the real environment, it is possible to use RL to learn a near-optimal policy from (usually a large amount of) simulation data. If the simulator is highly accurate, the learned policy should also perform well in the real environment.

Apart from some cases where the true environment and the simulator coincide (e.g., in game playing) or a nearly perfect simulator can be created from the law of physics (e.g., in simple control problems), in general we will need to construct the simulator using data from the real environment, making the overall approach an instance of *model-based RL*. As the algorithms for learning from simulated experience mature (which is what the RL community has mostly focused on), the bottleneck has shifted to the creation of a good simulator. *How can we learn a good model of the world from interaction experiences?*

A popular approach for meeting this challenge is to learn using a wide variety of simulators, which imparts robustness and adaptivity to the learned policies. Recent works have demonstrated the benefits of using such an ensemble of models, which can be used to either transfer policies from simulated to real-world domains, or to simply learn robust policies (Andrychowicz et al., 2018; Tobin et al., 2017; Rajeswaran et al., 2017). Borrowing the motivation from these empirical works, we notice that the process of learning a simulator inherently includes various choices like inductive biases, data collection policy, design aspects etc. As such, instead of relying on a sole approximate model for learning in simulation, interpolating between models obtained from different sources can provide better approximation of the real environment. Previous works like Buckman et al. (2018); Lee et al. (2019); Kurutach et al. (2018) have also demonstrated the effectiveness of using an ensemble of models for decreasing modelling error or its effect thereof during learning.

In this paper, we consider building an approximate model of the real environment from interaction data

using a set (or *ensemble*) of possibly inaccurate models, which we will refer to as the *base models*. The simplest way to combine the base models is to take a weighted combination, but such an approach is rather limited. For instance, each base model might be accurate in certain regions of the state space, in which case it is natural to consider a state-dependent mixture. We consider the problem of learning in such a setting, where one has to identify an appropriate combination of the base models through real-world interactions, so that the induced policy performs well in the real environment. The data collected through interaction with the real world can be a precious resource and, therefore, we need the learning procedure to be sample-efficient. Our main result is an algorithm that enjoys polynomial sample complexity guarantees, where the polynomial has no dependence on the size of the state and action spaces. We also study a more challenging setting where the featurization of states for learning the combination is unknown and has to be discovered from a large set of candidate features.

**Outline.** We formally set up the problem and notation in Section 2, and discuss related work in Section 3. The main algorithm is introduced in Section 4, together with its sample complexity guarantees. We then proceed to the feature selection problem in Section 5 and conclude in Section 6.

## 2 SETTING AND NOTATION

We consider episodic Markov decision processes (MDP). An MDP  $M$  is specified by a tuple  $(\mathcal{S}, \mathcal{A}, P, R, H, P_1)$ , where  $\mathcal{S}$  is the state space and  $\mathcal{A}$  is the action space.  $P$  denotes the transition kernel describing the system dynamics  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  and  $R$  is the per-timestep reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \Delta([0, 1])$ . The agent interacts with the environment for a fixed number of timesteps,  $H$ , which determines the horizon of the problem. The initial state distribution is  $P_1$ . The agent’s goal is to find a policy  $\pi : \mathcal{S} \times [H] \rightarrow \mathcal{A}$  which maximizes the value of the policy:

$$v_M^\pi := \mathbb{E}_{s \sim P_1} [V_{M,1}^\pi(s)]$$

where the value function at step  $h$  is defined as:

$$V_{M,h}^\pi(s) = \mathbb{E} \left[ \sum_{h'=h}^H r_{h'} \mid s_h = s, a_{h:H} \sim \pi, s_{h:H} \sim M \right]$$

Here we use “ $s_{h:H} \sim M$ ” to imply that the sequence of states are generated according to the dynamics of  $M$ . A policy is said to be optimal for  $M$  if it maximizes the value  $v_M^\pi$ . We denote such a policy as  $\pi_M$  and its value as  $v_M$ . We use  $M^*$  to denote the model of the

true environment, and use  $\pi^*$  and  $v^*$  as shorthand for  $\pi_{M^*}$  and  $v_{M^*}$ , respectively.

In our setting, the agent is given access to a set of  $K$  base MDPs  $\{M_1, \dots, M_K\}$ . They share the same  $\mathcal{S}, \mathcal{A}, H, P_1$ , and only differ in  $P$  and  $R$ . In addition, a feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{d-1}$  is given which maps state-action pairs to  $d$ -dimensional real vectors. Given these two objects, we consider the class of all models which can be obtained from the following state-dependent linear combination of the base models:

**Definition 2.1** (Linear Combination of Base Models). *For given model ensemble  $\{M_1, \dots, M_K\}$  and the feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{d-1}$ , we consider models parametrized by  $W$  with the following transition and reward functions:*

$$P^W(\cdot | s, a) = \sum_{k=1}^K [W\phi(s, a)]_k P^k(\cdot | s, a),$$

$$R^W(\cdot | s, a) = \sum_{k=1}^K [W\phi(s, a)]_k R^k(\cdot | s, a).$$

We will use  $M(W)$  to denote such a model for any parameter  $W \in \mathcal{W}$  with  $\mathcal{W}_0 \equiv \{W \in [0, 1]^{K \times d} : \sum_{i=1}^K W_{ij} = 1 \text{ for all } j \in [d]\}$ .

For now, let’s assume that there exists some  $W^*$  such that  $M^* = M(W^*)$ , i.e., the true environment can be captured by our model class; we will relax this assumption shortly.

To develop intuition, consider a simplified scenario where  $d = 1$  and  $\phi(s, a) \equiv 1$ . In this case, the matrix  $W$  becomes a  $K \times 1$  stochastic vector, and the true environment is approximated by a linear combination of the base models.

**Example 2.2** (Global convex combination of models). *If the base models are combined using a set of constant weights  $w \in \Delta_{K-1}$ , then this is a special case of Definition 2.1 where  $d = 1$  and each state’s feature vector is  $\phi(s, a) \equiv 1$ .*

In the more general case of  $d > 1$ , we allow the combination weights to be a linear transformation of the features, which are  $W\phi(s, a)$ , and hence obtain more flexibility in choosing different combination weights in different regions of the state-action space. A special case of this more general setting is when  $\phi$  corresponds to a partition of the state-action space into multiple groups, and the linear combination coefficients are constant within each group.

**Example 2.3** (State space partition). *Let  $\mathcal{S} \times \mathcal{A} = \bigcup_{i \in [d]} \mathcal{X}_i$  be a partition (i.e.,  $\{\mathcal{X}_i\}$  are disjoint). Let  $\phi_i(s, a) = \mathbf{1}(s, a \in \mathcal{X}_i)$  for all  $i \in [d]$  where  $\mathbf{1}[\cdot]$  is the indicator function. This  $\phi$  satisfies the condition that*

$\phi(s, a) \in \Delta_{d-1}$ , and when combined with a set of base models, forms a special case of Definition 2.1.

**Goal.** We consider the popular PAC learning objective: with probability at least  $1 - \delta$ , the algorithm should output a policy  $\pi$  with value  $v_{M^*}^\pi \geq v^* - \epsilon$  by collecting  $\text{poly}(d, K, H, 1/\epsilon, \log(1/\delta))$  episodes of data. Importantly, here the sample complexity is not allowed to depend on  $|\mathcal{S}|$  or  $|\mathcal{A}|$ . However, the assumption that  $M^*$  lies in the class of linear models can be limiting and, therefore, we will allow some approximation error in our setting as follows:

$$\theta := \min_{W \in \mathcal{W}} \sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left\| P^*(\cdot|s, a) - P^W(\cdot|s, a) \right\|_1 + \left\| R^*(\cdot|s, a) - R^W(\cdot|s, a) \right\|_1 \quad (1)$$

We denote the optimal parameter attaining this value by  $W^*$ . The case of  $\theta = 0$  represents the *realizable* setting where  $M^* = M(W^*)$  for some  $W^* \in \mathcal{W}$ . When  $\theta \neq 0$ , we cannot guarantee returning a policy with the value close to  $v^*$ , and will have to pay an additional penalty term proportional to the approximation error  $\theta$ , as is standard in RL theory.

**Further Notations** Let  $\pi_W$  be a shorthand for  $\pi_{M(W)}$ , the optimal policy in  $M(W)$ . When referring to value functions and state-action distributions, we will use the superscript to specify the policy and use the subscript to specify the MDP in which the policy is evaluated. For example, we will use  $V_{W',h}^W$  to denote the value of  $\pi_W$  (the optimal policy for model  $M(W')$ ) when evaluated in model  $M(W)$  starting from timestep  $h$ . The term  $d_{W',h}^W$  denotes the state-action distribution induced by policy  $\pi_W$  at timestep  $h$  in the MDP  $M(W')$ . Furthermore, we will write  $V_{M^*,h}^W$  and  $d_{M^*,h}^W$  when the evaluation environment is  $M^*$ . For conciseness,  $V_{W,h}$  and  $Q_{W,h}$  will denote the optimal (state- and Q-) value functions in model  $M(W)$  at step  $h$  (e.g.,  $V_{W,h}(s) \equiv V_{W,h}^W(s)$ ). The expected return of a policy  $\pi$  in model  $M(W)$  is defined as:

$$v_W^\pi = \mathbb{E}_{s \sim P_1} [V_{M(W),1}^\pi(s)]. \quad (2)$$

We assume that the total reward  $\sum_{h=1}^H r_h$  lies in  $[0, 1]$  almost surely in all MDPs of interest and under all policies. Further, whenever used, any value function at step  $H + 1$  (e.g.,  $V_{W,H+1}^\pi$ ) evaluates to 0 for any policy and any model.

### 3 RELATED WORK

**MDPs with low-rank transition matrices** Yang and Wang (2019b,a); Jin et al. (2019) have recently considered structured MDPs whose transition matrices

admit low-rank factorization, and the left matrix in the factorization are known to the learner as state-action features (corresponding to our  $\phi$ ). Their environmental assumption is a special case of ours, where the transition dynamics of each base model  $P^k(\cdot|s, a)$  is *independent* of  $s$  and  $a$ , i.e., each base MDP can be fully specified by a single density distribution over  $\mathcal{S}$ . This special case enjoys many nice properties, such as the value function of any policy is also linear in state-action features, and the linear value-function class is closed under the Bellman update operators, which are heavily exploited in their algorithms and analyses. In contrast, none of these properties hold under our more general setup, yet we are still able to provide sample efficiency guarantees. That said, we do note that the special case allows these recent works to obtain stronger results: their algorithms are both statistically and computationally efficient (ours is only statistically efficient), and some of these algorithms work without knowing the  $K$  base distributions.<sup>1</sup> We leave the problem of utilizing this linear structure in a regret minimization framework (Jin et al., 2019) for future work.

**Contextual MDPs** Abbasi-Yadkori and Neu (2014); Modi et al. (2018); Modi and Tewari (2019) consider a setting similar to our Example 2.2, except that the linear combination coefficients are visible to the learner and the base models are unknown. Therefore, despite the similarity in environmental assumptions, the learning objectives and the resulting sample complexities are significantly different (e.g., their guarantees depend on  $|\mathcal{S}|$  and  $|\mathcal{A}|$ ).

**Bellman rank** Jiang et al. (2017) have identified a structural parameter called Bellman rank for exploration under general value-function approximation, and devised an algorithm called OLIVE whose sample complexity is polynomial in the Bellman rank. A related notion is the witness rank (the model-based analog of Bellman rank) proposed by Sun et al. (2019). While our algorithm and analyses draw inspiration from these works, our setting does not obviously yield low Bellman rank or witness rank.<sup>2</sup> We will also provide a more detailed comparison to Sun et al. (2019), whose algorithm is most similar to ours among the existing works, in Section 4.

<sup>1</sup>In our setting, not knowing the base models immediately leads to hardness of learning, as it is equivalent to learning a general MDP without any prior knowledge even when  $d = K = 1$ . This requires  $\Omega(|\mathcal{S}||\mathcal{A}|)$  sample complexity (Azar et al., 2012), which is vacuous as we are interested in solving problems with arbitrarily large state and action spaces.

<sup>2</sup>In contrast, the low-rank MDPs considered by Yang and Wang (2019b,a); Jin et al. (2019) do admit low Bellman rank and low witness rank.

**Mixtures/ensembles of models** The closest work to our setting is the multiple model-based RL (MMRL) architecture proposed by Doya et al. (2002) where they also decompose a given domain as a convex combination of multiple models. However, instead of learning the combination coefficients for a given ensemble, their method trains the model ensemble and simultaneously learns a mixture weight for each *base model* as a function of state features. Their experiments demonstrate that each model specialized for different domains of the state space where the environment dynamics is predictable, thereby, providing a justification for using convex combination of models for simulation. Further, the idea of combining different models is inherently present in Bayesian learning methods where a posterior approximation of the real environment is iteratively refined using interaction data. For instance, Rajeswaran et al. (2017) introduce the EPOpt algorithm which uses an ensemble of simulated domains to learn robust and generalizable policies. During learning, they adapt the ensemble distribution (convex combination) over source domains using data from the target domain to progressively make it a better approximation. Similarly, Lee et al. (2019) combine a set of parameterized models by adaptively refining the mixture distribution over the latent parameter space. Here, we study a relatively simpler setting where a finite number of such base models are combined and give a frequentist sample complexity analysis for our method.

## 4 ALGORITHM AND MAIN RESULTS

In this section we introduce the main algorithm that learns a near-optimal policy in the aforementioned setup with a  $\text{poly}(d, K, H, 1/\epsilon, \log(1/\delta))$  sample complexity. We will first give the intuition behind the algorithm, and then present the formal sample complexity guarantees. Due to space constraints, we present the complete proof in the appendix. For simplicity, we will describe the intuition for the realizable case with  $\theta = 0$  ( $P^* \equiv P^{W^*}$ ). The pseudocode (Algorithm 1) and the results are, however, stated for the general case of  $\theta \neq 0$ .

At a high level, our algorithm proceeds in iterations  $t = 1, 2, \dots$ , and gradually refines a *version space*  $\mathcal{W}_t$  of plausible parameters. Our algorithm follows an *explore-or-terminate* template and in each iteration, either chooses to explore with a carefully chosen policy or terminates with a near-optimal policy. For exploration in the  $t$ -th iteration, we collect  $n$  trajectories  $\{(s_1^{(i)}, a_1^{(i)}, r_1^{(i)}, s_2^{(i)}, \dots, s_H^{(i)}, a_H^{(i)}, r_H^{(i)})\}_{i \in [n]}$  following some exploration policy  $\pi_t$  (Line 7). A key component of the algorithm is to extract knowledge

---

### Algorithm 1 PAC Algorithm for Linear Model Ensembles

---

- 1: **Input:**  $\{M_1, \dots, M_K\}, \epsilon, \delta, \phi(\cdot, \cdot), \mathcal{W}_0$
- 2: **for**  $t \rightarrow 1, 2, \dots$  **do**
- 3:   Compute *optimistic model*  $W_t$  and set  $\pi_t$  to  $\pi_{W_t}$

$$W_t \leftarrow \arg \max_{W \in \mathcal{W}_{t-1}} V_W$$

- 4:   Estimate the value of  $\pi_t$  using  $n_{\text{eval}}$  trajectories:

$$\hat{v}_t := \frac{1}{n_{\text{eval}}} \sum_{h=1}^H r_h^{(i)} \quad (3)$$

- 5:   **if**  $v_{W_t} - \hat{v}_t \leq 3\epsilon/4 + (3\sqrt{dK} + 1)H\theta$  **then**
- 6:     Terminate and output  $\pi_t$
- 7:     Collect  $n$  trajectories using  $\pi_t : a_h \sim \pi_t(s_h)$
- 8:     Estimate the matrix  $\hat{Z}_t$  and  $\hat{y}_t$  as

$$\hat{Z}_t := \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \bar{V}_{t,h} \left( s_h^{(i)}, a_h^{(i)} \right) \phi \left( s_h^{(i)}, a_h^{(i)} \right)^\top \quad (4)$$

$$\hat{y}_t := \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H r_{h+1}^{(i)} + V_{t,h+1} \left( s_{h+1}^{(i)} \right) \quad (5)$$

- 9:   Update the version space to  $\mathcal{W}_t$  as the set:

$$\left\{ W \in \mathcal{W}_{t-1} : \left| \hat{y}_t - \langle W, \hat{Z}_t \rangle \right| \leq \frac{\epsilon}{12\sqrt{dK}} + H\theta \right\} \quad (6)$$


---

about  $W^*$  from these trajectories. In particular, for every  $h$ , the bag of samples  $\{s_{h+1}^{(i)}\}_{i \in [n]}$  may be viewed as an unbiased draw from the following distribution

$$\frac{1}{n} \sum_{i=1}^n P^{W^*}(\cdot | s_h^{(i)}, a_h^{(i)}). \quad (7)$$

The situation for rewards is similar and will be omitted in the discussion. So in principle we could substitute  $W^*$  in Eq.(7) with any candidate  $W$ , and if the resulting distribution differs significantly from the real samples  $\{s_{h+1}^{(i)}\}_{h \in [H], i \in [n]}$ , we can assert that  $W \neq W^*$  and eliminate  $W$  from the version space. However, the state space can be arbitrarily large in our setting, and comparing state distributions directly can be intractable. Instead, we project the state distribution in Eq.(7) using a (non-stationary) discriminator function  $\{f_{t,h}\}_{h=1}^H$  (which will be chosen later) and consider the

following scalar property

$$\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \mathbb{E}_{\substack{r \sim R^{W^*}(\cdot | s_h^{(i)}, a_h^{(i)}) \\ s' \sim P^{W^*}(\cdot | s_h^{(i)}, a_h^{(i)})}} [r + f_{t,h+1}(s')], \quad (8)$$

which can be effectively estimated by

$$\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H (r_h^{(i)} + f_{t,h+1}(s_{h+1}^{(i)})). \quad (9)$$

Since we have projected states onto  $\mathbb{R}$ , Eq.(9) is the average of scalar random variables and enjoys state-space-independent concentration. Now, in order to test the validity of a parameter  $W$  in a given version space, we compare the estimate in eq. 9 with the prediction given by  $M(W)$ , which is:

$$\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \mathbb{E}_{\substack{r \sim R^W(\cdot | s_h^{(i)}, a_h^{(i)}) \\ s' \sim P^W(\cdot | s_h^{(i)}, a_h^{(i)})}} [r + f_{t,h+1}(s')]. \quad (10)$$

As we consider a linear model class, by using linearity of expectations, Eq.(10) may also be written as:

$$\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H [W \phi(s_h^{(i)}, a_h^{(i)})]^\top [\bar{V}_{t,h}(s_h^{(i)}, a_h^{(i)})] \quad (11)$$

$$= \left\langle W, \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \bar{V}_{t,h}(s_h^{(i)}, a_h^{(i)}) \phi(s_h^{(i)}, a_h^{(i)})^\top \right\rangle, \quad (12)$$

where  $\langle A, B \rangle$  denotes  $\text{Tr}(A^\top B)$  for any two matrices  $A$  and  $B$ . In eq. 12,  $\bar{V}_{t,h}$  is a function that maps  $(s, a)$  to a  $K$  dimensional vector with each entry being

$$[\bar{V}_{t,h}(s, a)]_k := \mathbb{E}_{\substack{r \sim R^k(\cdot | s, a) \\ s' \sim P^k(\cdot | s, a)}} [r + f_{t,h+1}(s')]. \quad (13)$$

The intuition behind Eq.(11) is that for each fixed state-action pair  $(s_h^{(i)}, a_h^{(i)})$ , the expectation in Eq.(8) can be computed by first taking expectation of  $r + f_{t,h+1}(s')$  over the reward and transition distributions of each of the  $K$  base models—which gives  $\bar{V}_h$ —and then aggregating the results using the combination coefficients. Rewriting Eq.(11) as Eq.(12), we see that Eq.(8) can also be viewed as a linear measurement of  $W^*$ , where the measurement matrix is again  $\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \bar{V}_h(s_h^{(i)}, a_h^{(i)}) \phi(s_h^{(i)}, a_h^{(i)})^\top$ . Therefore, by estimating this measurement matrix and the outcome (Eq.(9)), we obtain an approximate linear equality constraint over  $\mathcal{W}_{t-1}$  and can eliminate any candidate  $W$  that violates such constraints. By using a finite sample concentration bound over the inner product, we get a linear inequality constraint to update the version space (Eq.(6)).

The remaining concern is to choose the exploration policy  $\pi_t$  and the discriminator function  $\{f_{t,h}\}$  to ensure that the linear constraint induced in each iteration is significantly different from the previous ones and induces deep cuts in the version space. We guarantee this by choosing  $\pi_t := \pi_{W_t}$  and  $f_{t,h} := V_{W_t,h}$ <sup>3</sup>, where  $W_t$  is the *optimistic model* as computed in Line 3. That is,  $W_t$  predicts the highest optimal value among all candidate models in  $\mathcal{W}_{t-1}$ . Following a terminate-or-explore argument, we show that as long as  $\pi_t$  is suboptimal, the linear constraint induced by our choice of  $\pi_t$  and  $\{f_{t,h}\}$  will significantly reduce the volume of the version space, and the iteration complexity can be bounded as  $\text{poly}(d, K)$  by an ellipsoid argument similar to that of Jiang et al. (2017). Similarly, the sample size needed in each iteration only depends polynomially on  $d$  and  $K$  and incurs no dependence on  $|\mathcal{S}|$  or  $|\mathcal{A}|$ , as we have summarized high-dimensional objects such as  $f_{t,h}$  (function over states) using low-dimensional quantities such as  $\bar{V}_{t,h}$  (vector of length  $K$ ).

The bound on the number of iterations and the number of samples needed per iteration leads to the following sample complexity result:

**Theorem 4.1** (PAC bound for Alg. 1). *In Algorithm 1, if  $n_{\text{eval}} := \frac{32H^2}{\epsilon^2} \log \frac{4T}{\delta}$  and  $n = \frac{1800d^2KH^2}{\epsilon^2} \log \frac{8dKT}{\delta}$  where  $T = dK \log \frac{2\sqrt{2}KH}{\epsilon} / \log \frac{5}{3}$ , with probability at least  $1 - \delta$ , the algorithm terminates after using at most*

$$\tilde{O}\left(\frac{d^3K^2H^2}{\epsilon^2} \log \frac{dKH}{\delta}\right) \quad (14)$$

*trajectories and returns a policy  $\pi_T$  with a value  $v^T \geq v^* - \epsilon - (3\sqrt{dK} + 2)H\theta$ .*

By setting  $d$  and  $K$  to appropriate values, we obtain the following sample complexity bounds as corollaries:

**Corollary 4.2** (Sample complexity for partitions). *Since the state-action partitioning setting (Example 2.3) is subsumed by the general setup, the sample complexity is again:*

$$\tilde{O}\left(\frac{d^3K^2H^2}{\epsilon^2} \log \frac{dKH}{\delta}\right) \quad (15)$$

**Corollary 4.3** (Sample complexity for global convex combination). *When base models are combined without any dependence on state-action features (Example 2.2), the setting is special case of the general setup with  $d = 1$ . Thus, the sample complexity is:*

$$\tilde{O}\left(\frac{K^2H^2}{\epsilon^2} \log \frac{KH}{\delta}\right) \quad (16)$$

Our algorithm, therefore, satisfies the requirement of learning a near-optimal policy without any dependence

<sup>3</sup>We use the simplified notation  $V_{t,h}$  for  $V_{W_t,h}$ .

on the  $|\mathcal{S}|$  or  $|\mathcal{A}|$ . Moreover, we can also account for the approximation error  $\theta$  but also incur a cost of  $(3\sqrt{dK} + 1)H\theta$  in the performance guarantee of the final policy. As we use the projection of value functions through the linear model class, we do not model the complete dynamics of the environment. This leads to an additive loss of  $3\sqrt{dK}H\theta$  in value in addition to the best achievable value loss of  $2H\theta$  (see Corollary A.4 in the appendix).

**Comparison to OLIME (Sun et al., 2019)** Our Algorithm 1 shares some structural similarity with the OLIME algorithm proposed by Sun et al. (2019), but there are also several important differences. First of all, OLIME in each iteration will pick a time step and take uniformly random actions during data collection, and consequently incur polynomial dependence on  $|\mathcal{A}|$  in its sample complexity. In comparison, our main data collection step (Line 7) never takes a random deviation, and we do not pay any dependence on the cardinality of the action space. Secondly, similar to how we project the transition distributions onto a discriminator function (Eq.(7) and (8)), OLIME projects the distributions onto a *static discriminator class* and uses the corresponding integral probability metric (IPM) as a measure of model misfit. In our setting, however, we find that the most efficient and elegant way to extract knowledge from data is to use a *dynamic* discriminator function,  $V_{W_t, h}$ , which changes from iteration to iteration and depends on the previously collected data. Such a choice of discriminator function allows us to make direct cuts on the parameter space  $\mathcal{W}$ , whereas OLIME can only make cuts in the value prediction space.

**Computational Characteristics** In each iteration, our algorithm computes the optimistic policy within the version space. Therefore, we rely on access to the following *optimistic planning oracle*:

**Assumption 4.4** (Optimistic planning oracle). *We assume that when given a version space  $\mathcal{W}_t$ , we can obtain the optimistic model through a single oracle call for  $W_t = \arg \max_{W \in \mathcal{W}_t} V_W$ .*

It is important to note that any version space  $\mathcal{W}_t$  that we deal with is always an intersection of half-spaces induced by the linear inequality constraints. Therefore, one would hope to solve the optimistic planning problem in a computationally efficient manner given the nice geometrical form of the version space. However, even for a finite state-action space, we are not aware of any efficient solutions as the planning problem induces bilinear and non-convex constraints despite the linearity assumption. Many recently proposed algorithms also suffer from such a computational difficulty (Jiang et al., 2017; Dann et al., 2018; Sun et al., 2019).

Further, we also assume that for any given  $W$ , we can compute the optimal policy  $\pi_W$  and its value function: our elimination criteria in Eq. 6 uses estimates  $\hat{Z}_t$  and  $\hat{y}_t$  which in turn depend on the value function. This requirement corresponds to a standard planning oracle, and aligns with the motivation of our setting, as we can delegate these computations to any learning algorithm operating in the simulated environment with the given combination coefficient. Our algorithm, instead, focuses on careful and systematic exploration to minimize the sample complexity in the real world.

## 5 MODEL SELECTION WITH CANDIDATE PARTITIONS

In the previous section we showed that a near-optimal policy can be PAC-learned under our modeling assumptions, where the feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is given along with the approximation error  $\theta$ . In this section, we explore the more interesting and challenging setting where a realizable feature map  $\phi$  is unknown, but we know that the realizable  $\phi$  belongs to a candidate set  $\{\phi_i\}_{i=1}^N$ , i.e., the true environment satisfies our modeling assumption in Definition 2.1 under  $\phi = \phi_{i^*}$  for some  $i^* \in [N]$  with  $\theta_{i^*} = 0$ . Note that Definition 2.1 may be satisfied by multiple  $\phi_i$ 's; for example, adding redundant features to an already realizable  $\phi_{i^*}$  still yields a realizable feature map. In such cases, we consider  $\phi_{i^*}$  to be the most succinct feature map among all realizable ones, i.e., the one with the lowest dimensionality. Let  $d_i$  denote the dimensionality of  $\phi_i$ , and  $d^* = d_{i^*}$ .

One obvious baseline in this setup is to run Algorithm 1 with each  $\phi_i$  and select the best policy among the returned ones. This leads to a sample complexity of roughly  $\sum_{i=1}^N d_i^3$  (only the dependence on  $\{d_i\}_{i=1}^N$  is considered), which can be very inefficient: When there exists  $j$  such that  $d^* \ll d_j$ , we pay for  $d_j^3$  which is much greater than the sample complexity of  $d^*$ ; When  $\{d_i\}$  are relatively uniform, we pay a linear dependence on  $N$ , preventing us from competing with a large set of candidate feature maps.

So the key result we want to obtain is a sample complexity that scales as  $(d^*)^3$ , possibly with a mild multiplicative overhead dependence on  $d^*$  and/or  $N$  (e.g.,  $\log d^*$  and  $\log N$ ).

### Hardness Result for Unstructured $\{\phi_i\}$

Unfortunately, we show that this is impossible when  $\{\phi_i\}$  is unstructured via a lower bound. In the lower bound construction, we have an exponentially large set of candidate feature maps, all of which are state space partitions. Each of the partitions has trivial dimensionalities ( $d_i = 2, K = 2$ ), but the sample

complexity of learning is exponential, which can only be explained away as  $\Omega(N)$ .

**Proposition 5.1.** *For the aforementioned problem of learning an  $\epsilon$ -optimal policy using a candidate feature set of size  $N$ , no algorithm can achieve  $\text{poly}(d^*, K, H, 1/\epsilon, 1/\delta, N^{1-\alpha})$  sample complexity for any constant  $0 < \alpha < 1$ .*

On a separate note, besides providing formal justification for the structural assumption we will introduce later, this proposition is of independent interest as it also sheds light on the hardness of model selection with state abstractions. We discuss the further implications in Appendix B.

*Proof of Proposition 5.1.* We construct a linear class of MDPs with two base models  $M_1$  and  $M_2$  in the following way: Consider a complete tree of depth  $H$  with a branching factor of 2. The vertices forming the state space of  $M_1$  and  $M_2$  and the two outgoing edges in each state are the available actions. Both MDPs share the same deterministic transitions and each non-leaf node yields 0 reward. Every leaf node yields +1 reward in  $M_1$  and 0 in  $M_2$ . Now we construct a candidate partition set  $\{\phi_i\}$  of size  $2^H$ : for  $\phi_i$ , the  $i$ -th leaf node belongs to one equivalence class while all other leaf nodes belong to the other. (Non-leaf nodes can belong to either class as  $M_1$  and  $M_2$  agree on their transitions and rewards.)

Observe that the above model class contains a finite family of  $2^H$  MDPs, each of which only has 1 rewarding leaf node. Concretely, the MDP whose  $i$ -th leaf is rewarding is exactly realized under the feature map  $\phi_i$ , whose corresponding  $W^*$  is the identity matrix: the  $i$ -th leaf yields +1 reward as in  $M_1$ , and all other leaves yield 0 reward as in  $M_2$ . Learning in this family of  $2^H$  MDPs is provably hard (Krishnamurthy et al., 2016), as when the rewarding leaf is chosen adversarially, the learner has no choice but to visit almost all leaf nodes to identify the rewarding leaf as long as  $\epsilon$  is below a constant threshold. The proposition follows from the fact that in this setting  $d^* = 2$ ,  $K = 2$ ,  $1/\epsilon$  is a constant,  $N = 2^H$ , but the sample complexity is  $\Omega(2^H)$ .  $\square$

This lower bound shows the necessity of introducing structural assumptions in  $\{\phi_i\}$ . Below, we consider a particular structure of *nested partitions* that is natural and enables sample-efficient learning. Similar assumptions have also been considered in the state abstraction literature (e.g., Jiang et al., 2015).

### Nested Partitions as a Structural Assumption

Consider the case where every  $\phi_i$  is a partition. W.l.o.g. let  $d_1 \leq d_2 \leq \dots \leq d_N$ . We assume  $\{\phi_i\}$

is nested, meaning that  $\forall (s, a), (s', a')$ ,

$$\phi_i(s, a) = \phi_i(s', a') \implies \phi_j(s, a) = \phi_j(s', a'), \quad \forall i \leq j.$$

While this structural assumption almost allows us to develop sample-efficient algorithms, it is still insufficient as demonstrated by the following hardness result.

**Proposition 5.2.** *Fixing  $K = 2$ , there exist base models  $M_1$  and  $M_2$  and nested state space partitions  $\phi_1$  and  $\phi_2$ , such that it is information-theoretically impossible for any algorithm to obtain  $\text{poly}(d^*, H, K, 1/\epsilon, 1/\delta)$  sample complexity when an adversary chooses an MDP that satisfies our environmental assumption (Definition 2.1) under either  $\phi_1$  or  $\phi_2$ .*

*Proof.* We will again use an exponential tree style construction to prove the lower bound. Specifically, we construct two MDPs  $M$  and  $M'$  which are obtained by combining two base MDPs  $M_1$  and  $M_2$  using two different partitions  $\phi_1$  and  $\phi_2$ . The specification of  $M_1$  and  $M_2$  is exactly the same as in the proof of Proposition 5.1. We choose  $\phi_1$  to be a partition of size  $d_1 = 1$ , where all nodes are grouped together.  $\phi_2$  has size  $d_2 = 2^H$ , where each leaf node belongs to a separate group. (As before, which group the inner nodes belong to does not matter.)  $\phi_1$  and  $\phi_2$  are obviously nested. We construct  $M$  that is realizable under  $\phi_2$  by randomly choosing a leaf and setting the weights for the convex combination as  $(1/2 + 2\epsilon, 1/2 - 2\epsilon)$  for that leaf; for all other leaves, the weights are  $(1/2, 1/2)$ . This is equivalent to randomly choosing  $M$  from a set of  $2^H$  MDPs, each of which has only one *good* leaf node yielding a random reward drawn from  $\text{Ber}(1/2 + 2\epsilon)$  instead of  $\text{Ber}(1/2)$ . In contrast,  $M'$  is such that all leaf nodes yield  $\text{Ber}(1/2)$  reward, which is realizable under  $\phi_1$  with weights  $(1/2, 1/2)$ .

Observe that  $M$  and  $M'$  are exactly the same as the constructions in the proof of the multi-armed bandit lower bound by (Auer et al., 2002) (the number of arms is  $2^H$ ), where it has been shown that distinguishing between  $M$  and  $M'$  takes  $\Omega(2^H/\epsilon^2)$  samples. Now assume towards contradiction that there exists an algorithm that achieves  $\text{poly}(d^*, H, K, 1/\epsilon, 1/\delta)$  complexity; let  $f$  be the specific polynomial in its guarantee. After  $f(1, H, 2, 1/\epsilon, 1/\delta)$  trajectories are collected, the algorithm must stop if the true environment is  $M'$  to honor the sample complexity guarantee (since  $d^* = 1$ ,  $K = 2$ ), and proceed to collect more trajectories if  $M$  is the true environment (since  $d^* = 2^H$ ). Making this decision essentially requires distinguishing between  $M'$  and  $M$  using  $f(1, H, 2, 1/\epsilon, 1/\delta) = \text{poly}(H)$  trajectories, which contradicts the known hardness result from Auer et al. (2002). This proves the statement.  $\square$

Essentially, the lower bound creates a situation where

$d_1 \ll d_2$ , and the nature may adversarially choose a model such that either  $\phi_1$  or  $\phi_2$  is realizable. If  $\phi_1$  is realizable, the learner is only allowed a small sample budget and cannot fully explore with  $\phi_2$ , and if  $\phi_1$  is not realizable the learner must do the opposite. The information-theoretic lower bound shows that it is fundamentally hard to distinguish between the two situations: Once the learner explores with  $\phi_1$ , she cannot decide whether she should stop or move on to  $\phi_2$  without collecting a large amount of data.

This hardness result motivates our last assumption in this section, that the learner knows the value of  $v^*$  (a scalar) as side information. This way, the learner can compare the value of the returned policy in each round to  $v^*$  and effectively decide when to stop. This naturally leads to our Algorithm 2 that uses a doubling scheme over  $\{d_i\}$ , with the following sample complexity guarantee.

---

**Algorithm 2** Model Selection with Nested Partitions

---

**Input:**  $\{\phi_1, \phi_2, \dots, \phi_N\}, \{M_1, \dots, M_K\}, \epsilon, \delta, v^*$ .  
 $i \rightarrow 0$   
**while True do**  
 Choose  $\phi_i$  such that  $d_i$  is the largest among  $\{d_j : d_j \leq 2^i\}$ .  
 Run Algorithm 1 on  $\Phi_i$  with  $\epsilon_i = \frac{\epsilon}{2}$  and  $\delta_i = \frac{\delta}{2N}$ .  
 Terminate the sub-routine if  $t > d_i K \log \frac{2\sqrt{2KH}}{\epsilon} / \log \frac{5}{3}$ .  
 Let  $\pi_i$  be the returned policy (if any). Let  $\hat{v}_i$  be the estimated return of  $\pi_i$  using  $n_{\text{eval}} = \frac{9}{2\epsilon^2} \log \frac{2N}{\delta}$  Monte-Carlo trajectories.  
**if**  $\hat{v}_i \geq v^* - \frac{2\epsilon}{3}$  **then**  
 Terminate with output  $\pi_i$ .

---

**Theorem 5.3.** *When Algorithm 2 is run with the input  $v^*$ , with probability at least  $1 - \delta$ , it returns a near-optimal policy  $\pi$  with  $v^\pi \geq v^* - \epsilon$  using at most  $\tilde{O}\left(\frac{d^{*3}K^2H^2}{\epsilon^2} \log d^* \log \frac{d^*KHN}{\delta}\right)$  samples.*

*Proof.* In Algorithm 2, for each partition  $i$ , we run Algorithm 1 until termination or until the sample budget is exhausted. By union bound it is easy to verify that with probability at least  $1 - \delta$ , all calls to Algorithm 1 will succeed and the Monte-Carlo estimation of the returned policies will be  $(\epsilon/3)$ -accurate, and we will only consider this success event in the rest of the proof. When the partition under consideration is realizable, we get  $v_{M^*}^{\pi_i} \geq v^* - \epsilon/3$ , therefore

$$\hat{v}^i \geq v^{\pi_i} - \frac{\epsilon}{3} \geq v^* - \frac{2\epsilon}{3},$$

so the algorithm will terminate after considering a realizable  $\phi_i$ . Similarly, whenever the algorithm terminates, we have  $v^{\pi_i} \geq v^* - \epsilon$ . This is because

$$v^{\pi_i} \geq \hat{v}^i - \frac{\epsilon}{3} \geq v^* - \epsilon,$$

where the last inequality holds thanks to the termination condition of Algorithm 2, which relies on knowledge of  $v^*$ . The total number of iterations of the algorithm is at most  $\mathcal{O}(\log d^*)$ . Therefore, by taking a union bound over all possible iterations, the sample complexity is

$$\sum_{i=1}^J \tilde{O}\left(\frac{d_i^3 K^2 H^2}{\epsilon^2}\right) \leq \tilde{O}\left(\frac{d^{*3} K^2 H^2}{\epsilon^2} \log d^*\right). \quad \square$$

**Discussion.** Model selection in online learning—especially in the context of sequential decision making—is generally considered very challenging. There has been relatively limited work in the generic setting until recently for some special cases. For instance, Foster et al. (2019) consider the model selection problem in linear contextual bandits with a sequence of nested policy classes with dimensions  $d_1 < d_2 < \dots$ . They consider a similar goal of achieving sub-linear regret bounds which only scale with the optimal dimension  $d_{m^*}$ . In contrast to our result, they do not need to know the achievable value in the environment and give no-regret learning methods in the *knowledge-free* setting. However, this is not contradictory to our lower bound: Due to the extremely delayed reward signal, our construction is equivalent to a multi-armed bandit problem with  $2^H$  arms. Our negative result (Proposition 5.2) shows a lower bound on sample complexity which is exponential in horizon, therefore eliminating the possibility of sample efficient and knowledge-free model selection in MDPs.

## 6 CONCLUSION

In this paper, we proposed a sample efficient model based algorithms which learns a near-optimal policy by approximating the true environment via a feature dependent convex combination of a given ensemble. Our algorithm offers a sample complexity bound which is independent of the size of the environment and only depends on the number of parameters being learnt. In addition, we also consider a model selection problem, show exponential lower bounds and then give sample efficient methods under natural assumptions. The proposed algorithm and its analysis relies on a linearity assumption and shares this aspect with existing exploration methods for rich observation MDPs. We leave the possibility of considering a richer class of convex combinations to future work. Lastly, our work also revisits the open problem of coming up with a computational and sample efficient model based learning algorithm.



## Acknowledgements

This work was supported in part by a grant from the Open Philanthropy Project to the Center for Human-Compatible AI, and in part by NSF grant CAREER IIS-1452099.

## References

- Abbasi-Yadkori, Y. and Neu, G. (2014). Online learning in mdps with side information. *arXiv preprint arXiv:1406.6812*.
- Agarwal, A., Rakhlin, A., and Bartlett, P. (2008). Matrix regularization techniques for online multitask learning. Technical Report UCB/EECS-2008-138, EECS Department, University of California, Berkeley.
- Andrychowicz, M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., et al. (2018). Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256.
- Azar, M. G., Munos, R., and Kappen, H. J. (2012). On the sample complexity of reinforcement learning with a generative model. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 1707–1714. Omnipress.
- Buckman, J., Hafner, D., Tucker, G., Brevdo, E., and Lee, H. (2018). Sample-efficient reinforcement learning with stochastic ensemble value expansion. In *Advances in Neural Information Processing Systems*, pages 8224–8234.
- Christiano, P., Shah, Z., Mordatch, I., Schneider, J., Blackwell, T., Tobin, J., Abbeel, P., and Zaremba, W. (2016). Transfer from simulation to real world through learning deep inverse dynamics model. *arXiv preprint arXiv:1610.03518*.
- Dann, C., Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. (2018). On oracle-efficient pac rl with rich observations. In *Advances in Neural Information Processing Systems*, pages 1422–1432.
- Doya, K., Samejima, K., Katagiri, K.-i., and Kawato, M. (2002). Multiple model-based reinforcement learning. *Neural computation*, 14(6):1347–1369.
- Foster, D. J., Krishnamurthy, A., and Luo, H. (2019). Model selection for contextual bandits. *arXiv preprint arXiv:1906.00531*.
- Givan, R., Dean, T., and Greig, M. (2003). Equivalence notions and model minimization in markov decision processes. *Artificial Intelligence*, 147(1-2):163–223.
- Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. (2017). Contextual decision processes with low bellman rank are pac-learnable. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1704–1713. JMLR. org.
- Jiang, N., Kulesza, A., and Singh, S. (2015). Abstraction selection in model-based reinforcement learning. In *International Conference on Machine Learning*, pages 179–188.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2019). Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*.
- Kearns, M. and Singh, S. (2002). Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232.
- Krishnamurthy, A., Agarwal, A., and Langford, J. (2016). Pac reinforcement learning with rich observations. In *Advances in Neural Information Processing Systems 29*, pages 1840–1848.
- Kurutach, T., Clavera, I., Duan, Y., Tamar, A., and Abbeel, P. (2018). Model-ensemble trust-region policy optimization. In *International Conference on Learning Representations*.
- Lee, G., Hou, B., Mandalika, A., Lee, J., and Srinivasa, S. S. (2019). Bayesian policy optimization for model uncertainty. In *International Conference on Learning Representations*.
- Li, L., Walsh, T. J., and Littman, M. L. (2006). Towards a unified theory of state abstraction for mdps. In *ISAAC*.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2015). Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529.
- Modi, A., Jiang, N., Singh, S., and Tewari, A. (2018). Markov decision processes with continuous side information. In *Algorithmic Learning Theory*, pages 597–618.
- Modi, A. and Tewari, A. (2019). Contextual markov decision processes using generalized linear models. *arXiv preprint arXiv:1903.06187*.
- Mordatch, I., Mishra, N., Eppner, C., and Abbeel, P. (2016). Combining model-based policy search with online model learning for control of physical humanoids. In *2016 IEEE International Conference*

on *Robotics and Automation (ICRA)*, pages 242–248. IEEE.

- Rajeswaran, A., Ghotra, S., Ravindran, B., and Levine, S. (2017). Epopt: Learning robust neural network policies using model ensembles. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484.
- Sun, W., Jiang, N., Krishnamurthy, A., Agarwal, A., and Langford, J. (2019). Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. In Beygelzimer, A. and Hsu, D., editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 2898–2933, Phoenix, USA. PMLR.
- Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., and Abbeel, P. (2017). Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30. IEEE.
- Todorov, E., Erez, T., and Tassa, Y. (2012). Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE.
- Whitt, W. (1978). Approximations of dynamic programs, i. *Mathematics of Operations Research*, 3(3):231–243.
- Yang, L. and Wang, M. (2019a). Sample-optimal parametric q-learning using linearly additive features. In Chaudhuri, K. and Salakhutdinov, R., editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6995–7004, Long Beach, California, USA. PMLR.
- Yang, L. F. and Wang, M. (2019b). Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*.

## A Proofs from the main text

In this section, we provide a detailed proof as well as the key ideas used in the analysis. The proof uses an optimism based template which guarantees that either the algorithm terminates with a near-optimal policy or explores appropriately in the environment. We can show a polynomial sample complexity bound as the algorithm explores for a bounded number of iterations and the number of samples required in each iteration is polynomial in the desired parameters. We start with the key lemmas used in the analysis in Section A.1 with the final proof of the main theorem in Section A.2.

**Notation.** As in the main text, we use  $\langle X, Y \rangle$  for  $\text{Tr}(X^\top Y)$ . The notation  $\|A\|_F$  denotes the Frobenius norm  $\text{Tr}(A^\top A)$ . For any matrix  $A = (A^1 A^2 \dots A^n)$  in  $\mathbb{R}^{m \times n}$  with columns  $A^i \in \mathbb{R}^m$ , we will use  $\|A\|_{p,q}$  as the group norm:  $\|(\|A^1\|_p, \|A^2\|_p, \dots, \|A^n\|_p)\|_q$ .

### A.1 Key lemmas used in the analysis

For our analysis, we first define a term  $\mathcal{E}(W, h)$  for any parameter  $W$  which intuitively quantifies the model error at step  $h$ :

$$\mathcal{E}(W, h) := \mathbb{E}_{d_{M^*, h}^W} [\mathbb{E}_{M(W)} [r_h + V_{W, h+1}(s_{h+1}) | s_h, a_h] - \mathbb{E}_{M^*} [r_h + V_{W, h+1}(s_{h+1}) | s_h, a_h]] \quad (17)$$

We start with the following lemma which allows us to express the value loss by using a model  $M(W)$  in terms of these per-step quantities.

**Lemma A.1** (Value decomposition). *For any  $W \in \mathcal{W}$ , we can write the difference in two values:*

$$v_W - v_{M^*}^W = \mathcal{E}(W) := \sum_{h=1}^H \mathcal{E}(W, h) \quad (18)$$

*Proof.* We start with the value difference on the lhs:

$$\begin{aligned} v_W - v_{M^*}^W &= \mathbb{E}_{d_{M^*, 1}^W} [\mathbb{E}_{M(W)} [r_1 + V_{W, 2}(s_2) | s_1, a_1] - \mathbb{E}_{M^*} [r_1 + V_{M^*, 2}^W(s_2) | s_1, a_1]] \\ &= \mathbb{E}_{d_{M^*, 1}^W} [\mathbb{E}_{M(W)} [r_1 + V_{W, 2}(s_2) | s_1, a_1] - \mathbb{E}_{M^*} [r_1 + V_{W, 2}(s_2) - V_{W, 2}(s_2) + V_{M^*, 2}^W(s_2) | s_1, a_1]] \\ &= \mathbb{E}_{d_{M^*, 1}^W} [\mathbb{E}_{M(W)} [r_1 + V_{W, 2}(s_2) | s_1, a_1] - \mathbb{E}_{M^*} [r_1 + V_{W, 2}(s_2) | s_1, a_1]] \\ &\quad + \mathbb{E}_{d_{M^*, 2}^W} [V_{W, 2}(s_2) - V_{M^*, 2}^W(s_2)] \\ &= \mathcal{E}(W, 1) + \mathbb{E}_{d_{M^*, 2}^W} [\mathbb{E}_{M(W)} [r_2 + V_{W, 3}(s_3) | s_2, a_2] - \mathbb{E}_{M^*} [r_2 + V_{M^*, 3}^W(s_3) | s_2, a_2]] \end{aligned}$$

Unrolling the second expected value similarly till  $H$  leads to the desired result. □

At various places in our analysis, we will use the well-known simulation lemma to compare the value of a policy  $\pi$  across two MDPs:

**Lemma A.2** (Simulation Lemma (Kearns and Singh, 2002; Modi et al., 2018)). *Let  $M_1$  and  $M_2$  be two MDPs with the same state-action space. If the transition dynamics and reward functions of the two MDPs are such that:*

$$\begin{aligned} \|P^1(\cdot | s, a) - P^2(\cdot | s, a)\|_1 &\leq \epsilon_p \quad \forall s \in \mathcal{S}, a \in \mathcal{A} \\ |R^1(\cdot | s, a) - R^2(\cdot | s, a)| &\leq \epsilon_r \quad \forall s \in \mathcal{S}, a \in \mathcal{A} \end{aligned}$$

*then, for every policy  $\pi$ , we have:*

$$|v_{M_1}^\pi - v_{M_2}^\pi| \leq H\epsilon_p + \epsilon_r \quad (19)$$

Now, we will first use the assumption about linearity to prove the following key lemma of our analysis:

**Lemma A.3** (Decomposition of  $\mathcal{E}(W)$ ). *If  $\theta$  is the approximation error defined in eq. 1, then the quantity  $\mathcal{E}(W)$  can be bounded as follows:*

$$\mathcal{E}(W) \leq \langle W - W^*, \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^W} [\bar{V}_{W,h}(s_h, a_h) \phi(s_h, a_h)^\top] \rangle + H\theta \quad (20)$$

where  $\bar{V}_{W,h}(s_h, a_h) \in [0, 1]^K$  is a vector with the  $k^{\text{th}}$  entry as  $\mathbb{E}_{M_k} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]$ .

*Proof.* Using the definition of  $\mathcal{E}(W, h)$  from eq. 17, we rewrite the term as:

$$\begin{aligned} \mathcal{E}(W) &= \sum_{h=1}^H \mathcal{E}(W, h) \\ &= \mathbb{E}_{d_{M^*,h}^W} [\mathbb{E}_{M(W)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] - \mathbb{E}_{M^*} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]] \\ &= \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^W} [\mathbb{E}_{M(W)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] - \mathbb{E}_{M(W^*)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]] \\ &\quad + \mathbb{E}_{d_{M^*,h}^W} [\mathbb{E}_{M(W^*)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] - \mathbb{E}_{M^*} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]] \\ &\leq \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^W} [\mathbb{E}_{M(W)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] - \mathbb{E}_{M(W^*)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]] + \mathbb{E}_{d_{M^*,h}^W} [\theta] \end{aligned}$$

Here, we rewrite the inner expectation as:

$$\begin{aligned} \mathbb{E}_{M(W)} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] &= \sum_{k=1}^K (W \phi(s_h, a_h)) [k] \mathbb{E}_{M_k} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h] \\ &= \mathbb{E}_{d_{M^*,h}^W} [\langle W \phi(s_h, a_h), \bar{V}_{W,h}(s_h, a_h) \rangle] \\ &= \langle W, \mathbb{E}_{d_{M^*,h}^W} [\bar{V}_{W,h}(s_h, a_h) \phi(s_h, a_h)^\top] \rangle \end{aligned}$$

where  $\bar{V}_{W,h}(s_h, a_h) \in [0, 1]^K$  is a vector with the  $k^{\text{th}}$  entry as  $\mathbb{E}_{M_k} [r_h + V_{W,h+1}(s_{h+1}) | s_h, a_h]$ . Therefore, we can finally upper bound  $\mathcal{E}(W)$  by:

$$\mathcal{E}(W) \leq \langle W - W^*, \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^W} [\bar{V}_{W,h}(s_h, a_h) \phi(s_h, a_h)^\top] \rangle + H\theta$$

□

For conciseness, we use the notation  $V_{t,h}$  for the vector  $V_{W_t,h}$ . We write the matrix  $\sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^W} [\bar{V}_{W,h}(s_h, a_h) \phi(s_h, a_h)^\top]$  as  $Z_W$  and further use  $Z_t$  for  $Z_{W_t}$  which results in the bound:

$$\mathcal{E}(W) \leq \langle W - W^*, Z_W \rangle + H\theta$$

Further, using Lemma A.2, one can easily see the following result which we later use in Lemma A.5:

**Corollary A.4.** *For the true environment and the MDP  $M(W^*)$ , we have:*

$$\mathcal{E}(W^*) \leq |v_{W^*} - v_{M^*}^{W^*}| \leq H\theta \quad (21)$$

$$v_{M^*}^{W^*} \geq v^* - 2H\theta \quad (22)$$

*Proof.* Eq. 21 directly follows through from the assumption and Lemma A.2. For eq. 22, we have:

$$\begin{aligned} v_{M^*}^{W^*} &\geq v_{W^*} - H\theta \\ &\geq v_{W^*}^* - H\theta \\ &\geq v^* - 2H\theta \end{aligned}$$

□

By lemma A.1, we see that if the model-misfit error is controlled at each timestep, we can directly get a bound on the value loss incurred by using the greedy policy  $\pi_W$ . In Algorithm 1, we choose the optimistic policy  $W_t$  as the exploration policy which has the following property:

**Lemma A.5** (Explore-or-terminate). *If the estimate  $\hat{v}_t$  from eq. 3 satisfies the following inequality:*

$$\left| \hat{v}_t - v_{M^*}^{W_t} \right| \leq \frac{\epsilon}{4} \tag{23}$$

*throughout the execution of the algorithm and  $W^*$  is not eliminated from any  $W_t$  (version space is valid), then either of these two statements hold:*

- (i) *the algorithm terminates with output  $\pi_t$  such that  $v_{W^*}^{\pi_t} \geq v^* - (3\sqrt{dK} + 2)H\theta - \epsilon$*
- (ii) *the algorithm does not terminate and*

$$\mathcal{E}(W_t) \geq \frac{\epsilon}{2} + 3\sqrt{dK}H\theta + H\theta$$

*Proof.* If the algorithm doesn't terminate, then by the condition on line 5 and the assumption, we know that:

$$\mathcal{E}(W_t) = v_{W_t} - v_{W^*}^{W_t} \geq v_{W_t} - \hat{v}_t - \frac{\epsilon}{4} \geq \frac{\epsilon}{2} + 3\sqrt{dK}H\theta + H\theta$$

If the algorithm does terminate at step  $T$ , we have:

$$\begin{aligned} v_{M^*}^{\pi_T} &\geq \hat{v}_t - \epsilon/4 && \text{(Eq. 23)} \\ &\geq v_{W_T} - (3\sqrt{dK} + 1)H\theta - \epsilon && \text{(Alg. 1, Line 5)} \\ &\geq v_{W^*} - (3\sqrt{dK} + 1)H\theta - \epsilon && \text{(Optimism)} \\ &\geq v^* - (3\sqrt{dK} + 2)H\theta - \epsilon && \text{(Lemma A.4)} \end{aligned}$$

□

Lemma A.5 shows that either the algorithm terminates with a (near-)optimal policy or guarantees large model-misfit error for  $W_t$ . For bounding the number of iterations, we use a volumetric argument similar to Jiang et al. (2017). We will use the following Lemma to show the exponential rate of reduction in the volume of the version space:

**Lemma A.6** (Volume reduction for MVVEE, (Jiang et al., 2017)). *Consider a closed and bounded set  $V \subset \mathbb{R}^p$  and a vector  $a \in \mathbb{R}^p$ . Let  $B$  be any enclosing ellipsoid of  $V$  that is centered at the origin, and we abuse the same symbol for the symmetric positive definite matrix that defines the ellipsoid, i.e.,  $B = \{v \in \mathbb{R}^p : v^\top B^{-1}v \leq 1\}$ . Suppose there exists  $u \in V$  with  $|a^\top u| \geq \kappa$  and define  $B_+$  as the minimum volume enclosing ellipsoid of  $\{v \in B : |a^\top v| \leq \gamma\}$ . If  $\gamma/\kappa \leq 1/\sqrt{p}$ , we have*

$$\frac{\text{vol}(B_+)}{\text{vol}(B)} \leq \sqrt{p} \frac{\gamma}{\kappa} \left( \frac{p}{p-1} \right)^{(p-1)/2} \left( 1 - \frac{\gamma^2}{\kappa^2} \right)^{(p-1)/2} \tag{24}$$

Further, if  $\gamma/\kappa \leq \frac{1}{3\sqrt{p}}$ , the RHS of eq. 24 is less than 0.6.

In the following lemma, we now show that the exploration step can happen only a finite number of times:

**Lemma A.7** (Bounding the number of iterations). *If the estimates  $\hat{Z}_t$  and  $\hat{y}_t$  in eq. 4 and 5 satisfy:*

$$\left| \hat{y}_t - \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} \left[ \mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h] \right] \right| + \left| \langle W, \hat{Z}_t \rangle - \langle W, Z_t \rangle \right| \leq \frac{\epsilon}{12\sqrt{dK}} \tag{25}$$

*for all  $W \in \mathcal{W}$ , for all iterations in Algorithm 1, then  $W^*$  is never eliminated. Moreover, the number of exploration iterations of Algorithm 1 is at most  $T = dK \log \frac{2d\sqrt{KH}}{\epsilon} / \log \frac{5}{3}$ .*

*Proof.* By definition,  $W^* \in \mathcal{W}_0$ . We first show that  $W^*$  is never eliminated from the version space  $\mathcal{W}_t$ . Let  $\alpha_t := \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} [\mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h]]$  and by definition, we have  $\langle W^*, Z_t \rangle := \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} [\mathbb{E}_{M(W^*)} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h]]$ . Then, we have

$$\begin{aligned} |\hat{y}_t - \langle W^*, \hat{Z}_t \rangle| &\leq \left| \hat{y}_t - \alpha_t + \alpha_t - \langle W^*, Z_t \rangle + \langle W^*, Z_t \rangle - \langle W^*, \hat{Z}_t \rangle \right| \\ &\leq \frac{\epsilon}{12\sqrt{dK}} + \left| \alpha_t - \langle W^*, Z_t \rangle \right| \\ &\leq \frac{\epsilon}{12\sqrt{dK}} + H\theta \end{aligned}$$

Therefore,  $W^*$  always satisfies the update equation 6 and is never eliminated.

Now, we argue that the *volume* of the version space decreases at an exponential rate with each exploration iteration. To set up the volume reduction analysis, we first notice that:

$$\begin{aligned} \|W - W^*\|_F &\leq \sqrt{\sum_{i=1}^d \|W^i - W^{*i}\|_2^2} \\ &\leq \sqrt{2d} \end{aligned}$$

Therefore, the initial volume of a ball covering the space of flattened vectors in  $\mathcal{W}_0$  is at most  $c_{dK}(\sqrt{2d})^{dK}$ . We will now use Lemma A.6 by considering the flattened versions of the parameter matrices  $W$  in the dimension  $p = dK^4$ . Firstly, in each iteration until termination, we find a matrix  $W_t$  such that  $\mathcal{E}(W_t) \geq \frac{\epsilon}{2} + (3\sqrt{dK} + 1)H\theta$ . From Lemma A.3, we have:

$$\begin{aligned} \langle W_t - W^*, Z_t \rangle &\geq \mathcal{E}(W_t) - H\theta \\ &\geq \frac{\epsilon}{2} + 3\sqrt{dK}H\theta \end{aligned}$$

We will apply Lemma A.6 with  $W_t - W^*$  as the vector  $u$  and  $Z_t$  as the direction vector  $a$ . For the updated version space,  $\mathcal{W}_t$ , we have:

$$\begin{aligned} |\langle W - W^*, Z_t \rangle| &= |y_t - \hat{y}_t + \langle W, Z_t - \hat{Z}_t \rangle + \hat{Z}_t - \hat{y}_t| \\ &\leq \frac{\epsilon}{12\sqrt{dK}} + H\theta + \frac{\epsilon}{12\sqrt{dK}} \\ &\leq \frac{\epsilon}{6\sqrt{dK}} + H\theta \end{aligned}$$

Denoting  $B_{t-1}$  as the MVEE of the version space  $\mathcal{W}_{t-1}$ , we consider the MVEE  $B'_t$  of the set of vectors  $\mathcal{W}'_t \equiv \{W \in B_{t-1} : |\langle W - W^*, Z_t \rangle| \leq \epsilon/6\sqrt{dK} + H\theta\}$ . Clearly, we have  $\mathcal{W}_{t-1} \subseteq B_{t-1}$ , and hence,  $\mathcal{W}_t \subseteq \mathcal{W}'_t$ . By setting  $\kappa = \frac{\epsilon}{2}$  and  $\gamma = \frac{\epsilon}{3\sqrt{dK}}$  in Lemma A.6, we have:

$$\frac{\text{vol}(B_t)}{\text{vol}(B_{t-1})} \leq \frac{\text{vol}(B'_t)}{\text{vol}(B_{t-1})} \leq 0.6$$

This shows that the volume of the MVEE of the version spaces  $\mathcal{W}_t$  decreases with at least a constant rate. We now argue that the procedure stops after reaching a version space with sufficiently small volume.

For any  $W \in \mathcal{W}_t$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$  we have:

$$\begin{aligned} \left\| P^{W^*}(\cdot | s, a) - P^W(\cdot | s, a) \right\|_1 &\leq \left\| \sum_{k=1}^K [(W - W^*)\phi(s, a)]_k P^k(\cdot | s, a) \right\|_1 \\ &\leq \left\| (W - W^*)\phi(s, a) \right\|_1 \\ &\leq \sqrt{K} \left\| (W - W^*)\phi(s, a) \right\|_2 \\ &\leq \sqrt{dK} \|W - W^*\|_F \end{aligned}$$

<sup>4</sup>For avoiding ambiguity, we will directly use the matrix notations for inner products and norms.

Also, using lemma A.2, if the worst case error in next state transition estimates is bounded by  $\frac{\epsilon}{2H} + \theta$ , the optimistic value  $V_{W^*}^{\pi_W}$  is  $\epsilon + H\theta$ -optimal. Therefore, we only need to identify the matrix  $W$  to within  $\frac{\epsilon}{2H\sqrt{dK}}$  distance of  $W^*$ . Consequently, the terminating MVEE  $B_T$  satisfies:

$$B_T \supseteq \left\{ W : \|W - W^*\|_F \leq \frac{\epsilon}{2\sqrt{dKH}} \right\}$$

Therefore, we have  $\text{vol}(B_T) \geq c_{dK}(\epsilon/2\sqrt{dKH})^{dK}$ , and :

$$\frac{c_{dK}(\epsilon/2\sqrt{dKH})^{dK}}{c_{dK}(\sqrt{2d})^{dK}} \leq \frac{\text{vol}(B_T)}{\text{vol}(B_0)} \leq 0.6^T$$

By solving for  $T$ , we get that:

$$\begin{aligned} dK \log \frac{2\sqrt{2KH}}{\epsilon} &\geq T \log \frac{5}{3} \\ T &\leq dK \log \frac{2\sqrt{2KH}}{\epsilon} / \log \frac{5}{3} \end{aligned} \tag{26}$$

□

We now derive the number of trajectories required in each step to satisfy the validity requirements in Lemma A.5 and A.7:

**Lemma A.8** (Concentration for MC estimate  $\hat{v}_t$ ). *For any  $W_t \in \mathcal{W}$  with probability at least  $1 - \delta_1$ , we have:*

$$\left| \hat{v}_t - v_{M^*}^{W_t} \right| \leq \frac{\epsilon}{4} \tag{27}$$

if we set  $n_{eval} \geq \frac{8}{\epsilon^2} \log \frac{2}{\delta_1}$ .

*Proof.* Note that  $\hat{v}_t$  is an unbiased estimate of  $v_{M^*}^{W_t}$ . From our assumption on the expected sum of rewards, the return of each trajectory is bounded by 1 for all policies. Thus, the range of each summand for the estimate  $\hat{v}_t$  is  $[0, 1]$ . Then, the result follows from standard application of Hoeffding's inequality. □

**Lemma A.9** (Concentration of the model misfit error). *If  $n \geq \frac{1800d^2KH^2}{\epsilon^2} \log \frac{4dK}{\delta_2}$  in Algorithm. 1, then for a given  $t$ , each  $W \in \mathcal{W}$  and with probability at least  $1 - \delta_2$ , we have:*

$$\left| \hat{y}_t - \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} \left[ \mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h] \right] \right| + \left| \langle W, Z_t - \hat{Z}_t \rangle \right| \leq \frac{\epsilon}{12\sqrt{dK}} \tag{28}$$

*Proof.* To bound the first term, we note that  $\hat{y}_t$  is an unbiased estimate of  $\sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} \left[ \mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h] \right]$  and is bounded between  $[0, H]$ . Applying Hoeffding's inequality on the estimand  $\hat{y}_t$  when using  $n$  trajectories, with probability at least  $1 - \delta'$ , we get:

$$\left| \hat{y}_t - \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} \left[ \mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h] \right] \right| \leq H \sqrt{\frac{1}{2n} \log \frac{2}{\delta'}}$$

For bounding the second term, using Holder's inequality with matrix group norm (Agarwal et al., 2008), we first see:

$$\begin{aligned} |\langle W, \hat{Z}_t \rangle - \langle W, Z_t \rangle| &\leq \|W\|_{1,\infty} \|\hat{Z}_t - Z_t\|_{\infty,1} \\ &\leq \|\hat{Z}_t - Z_t\|_{\infty,1} \end{aligned}$$

We will now bound the estimation error for  $\widehat{Z}_t$ :

$$\begin{aligned} \|\widehat{Z}_t - Z_t\|_{\infty,1} &= \left\| \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \bar{V}_{t,h}(s_h^{(i)}, a_h^{(i)}) \phi(s_h^{(i)}, a_h^{(i)})^\top - \mathbb{E} \left[ \sum_{h=1}^H \bar{V}_{t,h}(s_h, a_h) \phi(s_h, a_h)^\top \right] \right\|_{\infty,1} \\ &= \sum_{j=1}^d \left\| \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \bar{V}_{t,h}(s_h^{(i)}, a_h^{(i)}) \phi(s_h^{(i)}, a_h^{(i)})[j] - \mathbb{E} \left[ \sum_{h=1}^H \bar{V}_{t,h}(s_h, a_h) \phi(s_h, a_h)[j] \right] \right\|_{\infty} \end{aligned}$$

We can consider each trajectory as a random sample from the distribution of trajectories induced by  $\pi_t$ . Therefore, by definition, each summand in the estimate of  $\widehat{Z}_t$  over  $n$  trajectories is an unbiased estimate of  $Z_t$ . Moreover, we know that each term in the entry  $\widehat{Z}_t[i, j]$  of the matrix  $\widehat{Z}_t$  is bounded by  $H$ . Using Bernstein's inequality for each term in the error matrix, and a union bound over all entries, with probability at least  $1 - \delta'$ , for all  $i, j$  we have:

$$|\widehat{Z}_t[i, j] - Z_t[i, j]| \leq \sqrt{\frac{2\text{Var}[\sum_{h=1}^H \bar{V}_{t,h}[i] \phi(s_h, a_h)[j]] \log \frac{2dK}{\delta'}}{n}} + \frac{2H \log \frac{2dK}{\delta'}}{n}$$

Summing up the maximum elements  $i_j$  of each column, we have:

$$\begin{aligned} \|\widehat{Z}_t - Z_t\|_{\infty,1} &\leq \sum_{j=1}^d \sqrt{\frac{2\text{Var}[\sum_{h=1}^H \bar{V}_{t,h}[i_j] \phi(s_h, a_h)[j]] \log \frac{2dK}{\delta'}}{n}} + \frac{2dH \log \frac{2dK}{\delta'}}{n} \\ &= \sum_{j=1}^d \sqrt{\frac{2\mathbb{E}[(\sum_{h=1}^H \phi(s_h, a_h)[j])^2] \log \frac{2dK}{\delta'}}{n}} + \frac{2dH \log \frac{2dK}{\delta'}}{n} \\ &\leq \sqrt{\frac{2d\mathbb{E}[\sum_{j=1}^d (\sum_{h=1}^H \phi(s_h, a_h)[j])^2] \log \frac{2dK}{\delta'}}{n}} + \frac{2dH \log \frac{2dK}{\delta'}}{n} \\ &\leq \sqrt{\frac{2d\mathbb{E}[(\sum_{j=1}^d \sum_{h=1}^H \phi(s_h, a_h)[j])^2] \log \frac{2dK}{\delta'}}{n}} + \frac{2dH \log \frac{2dK}{\delta'}}{n} \\ &\leq H \sqrt{\frac{2d \log \frac{2dK}{\delta'}}{n}} + \frac{2dH \log \frac{2dK}{\delta'}}{n} \end{aligned}$$

Here, for the first step, we have used the property that  $\bar{V}_{t,h}[i_j] \leq 1$ , the fact that variance is bounded by the second moment. The next step can be obtained by using Cauchy-Schwartz inequality. The last second step uses that property that for non-negative  $a_j$ ,  $\sum_j a_j^2 \leq (\sum_j a_j)^2$ . Now, if  $\frac{2d \log \frac{2}{\delta'}}{n} \leq 1$ , the above is bounded by  $2\sqrt{\frac{2d \log \frac{2}{\delta'}}{n}}$ .

Therefore, summing up the two terms with  $\delta' = \delta_2/2$ , with probability at least  $1 - \delta_2$  and for all  $W \in \mathcal{W}$ , we have:

$$\left| \hat{y}_t - \sum_{h=1}^H \mathbb{E}_{d_{M^*,h}^{W_t}} [\mathbb{E}_{M^*} [r_h + V_{t,h+1}(s_{h+1}) | s_h, a_h]] \right| + |\langle W, Z_t - \widehat{Z}_t \rangle| \leq \sqrt{\frac{8dH^2 \log \frac{4dK}{\delta_2}}{n}} + \sqrt{\frac{H^2 \log \frac{4}{\delta_2}}{2n}}$$

With some algebra, it can be verified that setting  $n = \frac{1800d^2KH^2}{\epsilon^2} \log \frac{4dK}{\delta_2}$  makes the total error bounded by  $\frac{\epsilon}{12\sqrt{dK}}$  with failure probability  $\delta_2$ .  $\square$

## A.2 Proof of Theorem 4.1

*Proof.* With the key lemmas in Section A.1, we can now prove the main result. For the main theorem, we need to ensure that the requirements in lemma A.5 and A.7 are satisfied. Since, our method maintains a version space of plausible weights, the validity of each iteration depends on every previous iteration being valid. Therefore, for a total failure probability of  $\delta$  for the algorithm, we assume:

- (i) Estimation of  $\hat{\mathcal{E}}(W_t)$  for all iterations: total failure probability  $\delta/2$



(ii) Updating the version space  $\mathcal{W}_t$ : total failure probability  $\delta/2$

We set  $\delta_1 = \delta/2T$  and  $\delta_2 = \delta/2T$  in Lemmas A.8 and A.9 respectively with  $T = dK \log \frac{2\sqrt{2KH}}{\epsilon} / \log \frac{5}{3}$ . By taking a union bound over maximum number of iterations, the total failure probability is bounded by  $\delta$ . Thus, the total number of trajectories unrolled by the algorithm is:

$$\begin{aligned} T(n_{\text{eval}} + n) &\leq \left( dK \log \frac{2\sqrt{2KH}}{\epsilon} / \log \frac{5}{3} \right) \left( \frac{8}{\epsilon^2} \log \frac{4T}{\delta} + \frac{1800d^2KH^2}{\epsilon^2} \log \frac{8dKT}{\delta} \right) \\ &= \tilde{O} \left( \frac{d^3K^2H^2}{\epsilon^2} \log \frac{dKH}{\delta} \right) \end{aligned}$$

Therefore, by the termination guarantee in Lemma A.5, we arrive at the desired upper bound on the number of trajectories required to guarantee a policy with value  $v_{M^*}^{\pi} \geq v^* - (3\sqrt{dK} + 1)H\theta$ .  $\square$

## B The Implication of Proposition 5.1 on the Hardness of Learning State Abstractions

Here we show that the proof of Proposition 5.1 can be adapted to show a related hardness result for learning with state abstractions. A state abstraction is a mapping  $\phi$  that maps the raw state space  $\mathcal{S}$  to some finite abstract state space  $\mathcal{S}_\phi$ , typically much smaller in size. When an abstraction  $\phi$  with good properties (e.g., preserving reward and transition dynamics) is known, one can leverage it in exploration and obtain a sample complexity that is polynomial in  $|\mathcal{S}_\phi|$  instead of  $|\mathcal{S}|$ . Among different types of abstractions, *bisimulation* (Whitt, 1978; Givan et al., 2003) is a very strict notion that comes with many nice properties (Li et al., 2006).

An open problem in state abstraction literature has been whether it is possible to perform model selection over a large set of candidate abstractions, i.e., designing an algorithm whose sample complexity only scales sublinearly (or ideally, logarithmically) with the cardinality of the candidate set. Using the construction from Proposition 5.1, we show that this is impossible without further assumptions:

**Proposition B.1.** *Consider a learner in an MDP that is equipped with a set of state abstractions,  $\{\phi_1, \phi_2, \dots, \phi_N\}$ . Each abstraction  $\phi_i$  maps the raw state space  $\mathcal{S}$  to a finite abstract state space  $\mathcal{S}_{\phi_i}$ . Even if there exists  $i^* \in [N]$  such that  $\phi^* = \phi_{i^*}$  is a bisimulation, no algorithm can achieve  $\text{poly}(|\mathcal{S}_{\phi^*}|, |\mathcal{A}|, H, 1/\epsilon, 1/\delta, N^{1-\alpha})$  sample complexity for any  $\alpha > 0$ .*

*Proof.* Following the proof of Proposition 5.1, we consider the family of MDPs that share the same deterministic transition dynamics with a complete tree structure, where each MDP only has one rewarding leaf. Let  $N = 2^H$  be the number of leaves, and let the MDPs be  $\{M_i\}$  where the index indicates the rewarding leaf. We then construct a set of  $N$  abstractions where one of them will always be a bisimulation regardless of which MDP we choose from the family. Consider the  $i$ -th abstraction,  $\phi_i$ . At each level  $h$ ,  $\phi_i$  aggregates the state on the optimal path in its own equivalence class, and aggregates all other states together. It does not aggregate states across levels. It is easy to verify that  $\phi_i$  is a bisimulation in MDP  $M_i$ , and  $|\mathcal{S}_{\phi_i}| = 2H$ . If the hypothetical algorithm existed, it would achieve a sample complexity sublinear in  $N$ , as  $N = 2^H$  and all other relevant parameters (e.g.,  $|\mathcal{S}_{\phi_i}|$  and  $H$ ) are at most polynomial in  $H$ . However, the set of abstractions  $\{\phi_i\}$  is uninformative in this specific problem and does not affect the  $\Omega(2^H)$  sample complexity, which completes the proof.  $\square$