# FLAMBE: Structural Complexity and Representation Learning of Low Rank MDPs

Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun*

Microsoft Research

**Abstract**

In order to deal with the curse of dimensionality in reinforcement learning (RL), it is common practice to make parametric assumptions where values or policies are functions of some low dimensional feature space. This work focuses on the representation learning question: how can we *learn* such features? Under the assumption that the underlying (unknown) dynamics correspond to a low rank transition matrix, we show how the representation learning question is related to a particular *non-linear* matrix decomposition problem. Structurally, we make precise connections between these low rank MDPs and latent variable models, showing how they significantly generalize prior formulations for representation learning in RL. Algorithmically, we develop FLAMBE, which engages in exploration and representation learning for provably efficient RL in low rank transition models.

## 1 Introduction

The ability to learn effective transformations of complex data sources, sometimes called representation learning, is an essential primitive in modern machine learning, leading to remarkable achievements in language modeling, vision, and serving as a partial explanation for the success of deep learning more broadly (Bengio et al., 2013). In Reinforcement Learning (RL), several works have shown empirically that learning succinct representations of perceptual inputs can accelerate the search for decision-making policies (Pathak et al., 2017; Tang et al., 2017; Oord et al., 2018; Srinivas et al., 2020). However, representation learning for RL is far more subtle than it is for supervised learning (Du et al., 2019a; Van Roy and Dong, 2019; Lattimore and Szepesvari, 2019), and the theoretical foundations of representation learning for RL are nascent.

The first question that arises in this context is: what is a good representation? Intuitively, a good representation should help us achieve greater sample efficiency on downstream tasks. For supervised learning, several theoretical works adopt the perspective that a good representation should permit simple models to achieve high accuracy on tasks of interest (Baxter, 2000; Maurer et al., 2016; Arora et al., 2019; Tosh et al., 2020). Lifting this perspective to reinforcement learning, it is natural to ask that we can express value functions and policies as simple functions of our representation. This may allow us to leverage recent work on sample efficient RL with parametric function approximation.

The second question is: how do we *learn* such a representation when it is not provided in advance? This question is particularly challenging because representation learning is intimately tied to exploration. We cannot learn a good representation without a comprehensive dataset of experience from the environment, but a good representation may be critical for efficient exploration.

This work considers these questions in the context of low rank MDPs (Jiang et al., 2017) (also known as factorizing MDPs (Rendle et al., 2010), factored linear MDPs (Yao et al., 2014), and linear MDPs (Jin

---

*alekha@microsoft.com, sham@cs.washington.edu, akshaykr@microsoft.com, sun.wen@microsoft.com

| Algorithm | Setting | Sample Complexity | Computation |
|---|---|---|---|
| PCID (Du et al., 2019b) | block MDP | $d^4 H^2 K^4 \left( \frac{1}{\eta^4 \gamma^2} + \frac{1}{\varepsilon^2} \right)$ | Oracle efficient |
| HOMER (Misra et al., 2019) | block MDP | $d^8 H^4 K^4 \left( \frac{1}{\eta^3} + \frac{1}{\varepsilon^2} \right)$ | Oracle efficient |
| OLIVE (Jiang et al., 2017) | low Bellman rank | $\frac{d^2 H^3 K}{\varepsilon^2}$ | Inefficient |
| Sun et al. (2019) | low Witness rank | $\frac{d^2 H^3 K}{\varepsilon^2}$ | Inefficient |
| FLAMBE (this paper) | low rank MDP | $\frac{d^7 K^9 H^{22}}{\varepsilon^{10}}$ | Oracle efficient |

Table 1: Comparison of methods for representation learning in RL. Settings from least to most general are: block MDP, low rank MDP, low Bellman rank, low Witness rank. In all cases $d$ is the embedding dimension, $H$ is the horizon, $K$ is the number of actions, $\eta$ and $\gamma$ parameterize reachability and margin assumptions, and $\varepsilon$ is the accuracy. Dependence on function classes and logarithmic factors are suppressed. Oracle and realizability assumptions vary. Block MDP algorithms discover a *one-hot* representation to discrete latent states. Bellman/Witness rank approaches can take a class $\Phi$ of embedding functions and search over simple policies or value functions composed with $\Phi$ (see Section 4 and Appendix A.3 for details).

et al., 2019; Yang and Wang, 2019b)), which we argue provide a natural framework for studying representation learning in RL. Concretely, these models assume there exists low dimensional embedding functions $\phi(x, a), \mu(x')$ such that the transition operator $T$ satisfies $T(x' \mid x, a) = \langle \phi(x, a), \mu(x') \rangle$, where $T(x' \mid x, a)$ specifies the probability of the next state $x'$ given the previous state $x$ and action $a$. Low rank MDPs address the first issue above (on what constitutes a good representation) in that if the features $\phi$ are known to the learner, then sample efficient learning is possible (Jin et al., 2019; Yang and Wang, 2019b).

**Our contributions.** We address the question of learning the representation $\phi$ in a low rank MDP. To this end our contributions are both structural and algorithmic.

1. **Expressiveness of low rank MDPs.** Our algorithmic development leverages a re-formulation of the low rank dynamics in terms of an equally expressive, but more interpretable latent variable model. We provide several structural results for low rank MDPs, relating it to other models studied in prior work on representation learning for RL. In particular, we show that low rank MDPs are significantly more expressive than the block MDP model (Du et al., 2019b; Misra et al., 2019).

2. **Feature learning.** We develop a new algorithm, called FLAMBE for "Feature learning and model based exploration", that learns a representation for low rank MDPs. We prove that under realizability assumptions, FLAMBE learns a *uniformly accurate* model of the environment as well as a feature map that enables the use of linear methods for RL, in a statistically and computationally efficient manner. These guarantees enable downstream reward maximization, for *any* reward function, with no additional data collection.

Our results and techniques provide new insights on representation learning for RL and also significantly increase the scope for provably efficient RL with rich observations (see Table 1).

## 2 Low Rank MDPs

We consider an episodic Markov decision process $\mathcal{M}$ with episode length $H \in \mathbb{N}$, state space $\mathcal{X}$, and a finite action space $\mathcal{A} = \{1, \ldots, K\}$. In each episode, a trajectory $\tau = (x_0, a_0, x_1, a_1, \ldots, x_{H-1}, a_{H-1}, x_H)$ is generated, where (a) $x_0$ is a starting state, and (b) $x_{h+1} \sim T_h(\cdot \mid x_h, a_h)$, and (c) all actions $a_{0:H-1}$ are chosen by the agent. We assume the starting state is fixed and that there is only one available action at time $0$.[1] The operators $T_h : \mathcal{X} \times \mathcal{A} \to \Delta(\mathcal{X})$ denote the (non-stationary) transition dynamics for each time step.

---

[1]This easily accommodates the standard formulation with a non-degenerate initial distribution by defining $T_0(\cdot \mid x_0, a_0)$ to be the initial distribution. This setup is notationally more convenient, since we do not need special notation for the starting distribution.

As is standard in the literature, a policy $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ is a (randomized) mapping from states to actions. We use the notation $\mathbb{E}\left[\cdot \mid \pi, \mathcal{M}\right]$ to denote expectations over states and actions observed when executing policy $\pi$ in MDP $\mathcal{M}$. We abuse notation slightly and use $[H]$ to denote $\{0, \dots, H-1\}$.

**Definition 1.** *An operator $T : \mathcal{X} \times \mathcal{A} \to \Delta(\mathcal{X})$ admits a* low rank decomposition *with dimension $d \in \mathbb{N}$ if there exists two embedding functions $\phi^\star : \mathcal{X} \times \mathcal{A} \to \mathbb{R}^d$ and $\mu^\star : \mathcal{X} \to \mathbb{R}^d$ such that*

$$\forall x, x' \in \mathcal{X}, a \in \mathcal{A} : T(x' \mid x, a) = \langle \phi^\star(x,a), \mu^\star(x') \rangle .$$

*For normalization,[2] we assume that $\|\phi^\star(x,a)\|_2 \leq 1$ for all $x, a$ and for any function $g : \mathcal{X} \to [0,1]$, $\left\|\int \mu^\star(x) g(x) dx\right\|_2 \leq \sqrt{d}$. An MDP $\mathcal{M}$ is a* low rank MDP *if for each $h \in [H]$, $T_h$ admits a low rank decomposition with dimension $d$. We use $\phi_h^\star, \mu_h^\star$ to denote the embeddings for $T_h$.*

Throughout we assume that $\mathcal{M}$ is a low rank MDP with dimension $d$. Note that the condition on $\mu^\star$ ensures that the Bellman backup operator is well-behaved.

**Function approximation for representation learning.** We consider state spaces $\mathcal{X}$ that are arbitrarily large, so that some form of function approximation is necessary to generalize across states. For representation learning, it is natural to grant the agent access to two function classes $\Phi \subset \mathcal{X} \times \mathcal{A} \to \mathbb{R}^d$ and $\Upsilon \subset \mathcal{X} \to \mathbb{R}^d$ of candidate embeddings, which we can use to identify the true embeddings $(\phi^\star, \mu^\star)$. To facilitate this model selection task, we posit a *realizability* assumption.

**Assumption 1** (Realizability). *We assume that for each $h \in [H]$: $\phi_h^\star \in \Phi$ and $\mu_h^\star \in \Upsilon$.*

We desire sample complexity bounds that scale logarithmically with the cardinality of the classes $\Phi$ and $\Upsilon$, which we assume to be finite. Extensions that permit infinite classes with bounded statistical complexity (e.g., VC-classes) are not difficult.

In Appendix A, we show that the low rank assumption alone, without Assumption 1, is not sufficient for obtaining performance guarantees that are independent of the size of the state space. Hence, additional modeling assumptions are required, and we encode these in $\Phi, \Upsilon$.

**Learning goal.** We focus on the problem of reward-free exploration (Hazan et al., 2019; Jin et al., 2020), where the agent interacts with the environment with no reward signal. When considering model-based algorithms, a natural reward-free goal is *system identification*: given function classes $\Phi, \Upsilon$, the algorithm should learn a model $\widehat{\mathcal{M}} := (\hat{\phi}_{0:H-1}, \hat{\mu}_{0:H-1})$ that uniformly approximates the environment $\mathcal{M}$. We formalize this with the following performance criteria:

$$\forall \pi, h \in [H] : \mathbb{E}\left[\left\|\langle \hat{\phi}_h(x_h, a_h) \hat{\mu}_h(\cdot) \rangle - T_h(\cdot \mid x_h, a_h)\right\|_{\mathrm{TV}} \mid \pi, \mathcal{M}\right] \leq \varepsilon. \tag{1}$$

Here, we ask that our model accurately approximates the one-step dynamics from the state-action distribution induced by following *any* policy $\pi$ for $h$ steps in the real environment.

System identification also implies a quantitative guarantee on the learned representation $\hat{\phi}_{0:H-1}$: we can approximate the Bellman backup of any value function on any data-distribution.

**Lemma 1.** *If $\widehat{\mathcal{M}} = (\hat{\phi}_{0:H-1}, \hat{\mu}_{0:H-1})$ satisfies (1), then*

$$\forall h \in [H], V : \mathcal{X} \to [0,1], \exists \theta_h : \max_\pi \mathbb{E}\left[\left|\langle \theta_h, \hat{\phi}_h(x_h, a_h) \rangle - \mathbb{E}\left[V(x_{h+1}) \mid x_h, a_h\right]\right| \mid \pi, \mathcal{M}\right] \leq \varepsilon.$$

Thus, linear function approximation using our learned features suffices to fit the $Q$ function associated with any policy and explicitly given reward.[3] The guarantee also enables dynamic programming techniques for policy optimization. In other words, (1) verifies that we have found a good representation, in a quantitative sense, and enables tractable reward maximization for any known reward function.

---

[2] See the proof of Lemma B.1 in Jin et al. (2019) for this form of the normalization assumption.

[3] Formally, we append the immediate reward to the features.

# 3 Related work

Low rank models are prevalent in dynamics and controls (Thon and Jaeger, 2015; Littman and Sutton, 2002; Singh et al., 2004). The low rank MDP in particular has been studied in several works in the context of planning (Barreto et al., 2011; Barreto and Fragoso, 2011), estimation (Duan et al., 2020), and in the generative model setting (Yang and Wang, 2019a). Regarding nomenclature, to our knowledge the name *low rank MDP* appears first in Jiang et al. (2017), although Rendle et al. (2010) refer to it as *factorizing MDP*, Yao et al. (2014) call it a factored linear MDP, and Barreto et al. (2011) refer to a similar model as *stochastic factorization*. More recently, it has been called the *linear MDP* by Jin et al. (2019). We use *low rank MDP* because it highlights the key structural property of the dynamics, and because we study the setting where the embeddings are unknown, which necessitates non-linear function approximation.

Turning to reinforcement learning with function approximation and exploration, a large body of effort focuses on (essentially) linear methods (Yang and Wang, 2019b; Jin et al., 2019; Cai et al., 2019; Modi et al., 2020; Du et al., 2019c; Wang et al., 2019). Closest to our work are the results of Jin et al. (2019) and Yang and Wang (2019b), who consider low rank MDPs with known feature maps $\phi_{0:H-1}^{\star}$ (Yang and Wang (2019b) also assumes that $\mu_{0:H-1}^{\star}$ is known up to a linear map). These results are encouraging and motivate our representation learning formulation, but, on their own, these methods cannot leverage the inductive biases provided by neural networks to scale to rich state spaces.

There are methods for more general, non-linear, function approximation, but these works either (a) require strong environment assumptions such as determinism (Wen and Van Roy, 2013; Du et al., 2020), (b) require strong function class assumptions such as bounded Eluder dimension (Russo and Van Roy, 2013; Osband and Van Roy, 2014), (c) have sample complexity scaling linearly with the function class size (Lattimore et al., 2013; Ortner et al., 2014) or (d) are computationally intractable (Jiang et al., 2017; Sun et al., 2019; Dong et al., 2019). Note that Ortner et al. (2014); Jiang et al. (2015) consider a form of representation learning, abstraction selection, but the former scales linearly with the number of candidate abstractions, while the latter does not address exploration.

**Bellman/Witness rank.** We briefly expand on this final category of computationally inefficient methods. For model-free reinforcement learning, Jiang et al. (2017) give an algebraic condition, in terms of a notion called the Bellman rank, on the environment and a given function approximation class, under which sample efficient reinforcement learning is always possible. Sun et al. (2019) extend the definition to model-based approaches, with the notion of Witness rank. As we will see in the next section, the low rank MDP with a function class derived from $\Phi$ (and $\Upsilon$) admits low Bellman (resp., Witness) rank, and so these results imply that our setting is statistically tractable.

**Block MDPs.** Finally, we turn to theoretical works on representation learning for RL. Du et al. (2019b) introduce the *block MDP* model, in which there is a finite latent state space $\mathcal{S}$ that governs the transition dynamics, and each "observation" $x \in \mathcal{X}$ is associated with a latent state $s \in \mathcal{S}$, so the state is *decodable*. The natural representation learning goal is to recover the latent states, and Du et al. (2019b); Misra et al. (2019) show that this can be done, in concert with exploration, in a statistically and computationally efficient manner. Since the block MDP can be easily expressed as a low rank MDP, our results can be specialized to this setting, where they yield comparable guarantees. On the other hand, we will see that the low rank MDP is significantly more expressive, and so our results greatly expand the scope for provably efficient representation learning and reinforcement learning.

# 4 Expressiveness of low rank MDPs

Before turning to our algorithmic development, we discuss connections between low-rank MDPs and models studied in prior work. This discussion is facilitated by formalizing a connection between MDP transition operators and latent variable graphical models.
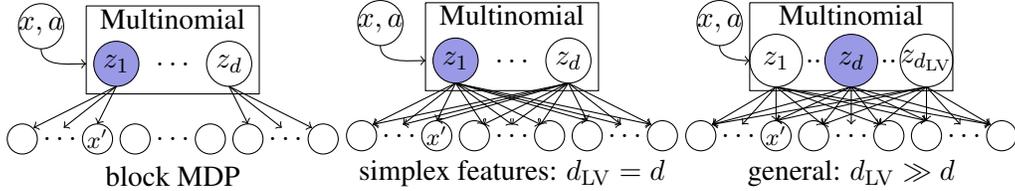
Figure 2: The latent variable interpretation of low rank MDPs, where $(x, a)$ induces a distribution over latent variable $z$. Left: in block MDPs, latent variables induce a partition over the next state $x'$. Center: simplex features have embedding dimension equal to the number of latent variables. Right: low rank MDPs can have exponentially more latent variables than the dimension, $d_{\mathrm{LV}} \gg d$.

**Definition 2.** *The* latent variable representation *of a transition operator* $T : \mathcal{X} \times \mathcal{A} \to \Delta(\mathcal{X})$ *is a latent space* $\mathcal{Z}$ *along with functions* $\psi : \mathcal{X} \times \mathcal{A} \to \Delta(\mathcal{Z})$ *and* $\nu : \mathcal{Z} \to \Delta(\mathcal{X})$, *such that* $T(\cdot \mid x, a) = \int \nu(\cdot \mid z)\psi(z \mid x, a)dz$. *The* latent variable dimension *of* $T$, *denoted* $d_{\mathrm{LV}}$ *is the cardinality of smallest latent space* $\mathcal{Z}$ *for which* $T$ *admits a latent variable representation.*

See Figure 1. In this representation, (1) each $(x, a)$ pair induces a "posterior" distribution $\psi(x, a) \in \Delta(\mathcal{Z})$ over $z$, (2) we sample $z \sim \psi(x, a)$, and (3) then sample $x' \sim \nu(\cdot \mid z)$, where $\nu$ specifies the "emission" distributions. As notation, we typically write $\nu(x) \in \mathbb{R}^{\mathcal{Z}}$ with coordinates $\nu(x)[z] = \nu(x \mid z)$ and we call $\psi, \nu$ the *simplex features*, following the example described by Jin et al. (2019). When considering $H$-step MDPs, this representation allows us to augment the trajectory $\tau$ with the latent variables $\tau = (x_0, a_0, z_1, x_1, \ldots z_{H-1}, x_{H-1}, a_{H-1}, x_H)$. Here note that $z_h$ is the latent variable that generates $x_h$.
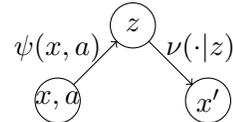


Figure 1: The latent variable interpretation.

Note that all transition operators admit a trivial latent variable representation, as we may always take $\psi(x, a) = T(\cdot \mid x, a)$. However, when $T$ is endowed with additional structure, the latent variable representations are more interesting. For example, this viewpoint already certifies a factorization $T(x' \mid x, a) = \langle \psi(x, a), \nu(x') \rangle$ with embedding dimension $|\mathcal{Z}|$, and so $d_{\mathrm{LV}}$ (if it is finite) is an upper bound on the rank of the transition operator. On the other hand, compared with Definition 1, this factorization additionally requires that $\psi(x, a)$ and $\nu(\cdot \mid z)$ are probability distributions. Since the factorization is non-negative, $d_{\mathrm{LV}}$ is the *non-negative rank* of the transition operator.

The latent variable representation enables a natural comparison of the expressiveness of various models, and, as we will see in the next section, yields insights that facilitate algorithm design. We now examine models that have been introduced in prior works and their properties relative to Definition 1.

**Block MDPs.** A block MDP (Du et al., 2019b; Misra et al., 2019) is clearly a latent variable model with $\mathcal{Z}$ corresponding to the latent state space $\mathcal{S}$ and the additional restriction that two latent variables $z$ and $z'$ have disjoint supports in their respective emissions $\nu(\cdot \mid z)$ and $\nu(\cdot \mid z')$ (see the left panel of Figure 2). Therefore, a block MDP is a low rank MDP with rank $d \leq |\mathcal{S}|$, but the next result shows that a low rank MDP is significantly more expressive.

**Proposition 1.** *For any* $d \geq 2$ *and any* $M \in \mathbb{N}$ *there exists an environment on* $|\mathcal{X}| = M$ *states, that can be expressed as a low rank MDP with embedding dimension* $d$, *but for which any block MDP representation must have* $M$ *latent states.*

In fact, the MDP that we construct for the proof, admits a latent variable representation with $|\mathcal{Z}| = d$, but does not admit a non-trivial block MDP representation. This separation exploits the decodability restriction of block MDPs, which is indeed quite limiting in terms of expressiveness.

5

**Simplex features.** Given the latent variable representation and the fact that it certifies a rank of at most $d_{\mathrm{LV}}$, it is natural to ask if this representation is canonical for all low rank MDPs. In other words, for any transition operator with rank $d$, can we express it as a latent variable model with $|\mathcal{Z}| = d$, or equivalently with simplex features of dimension $d$?

As discussed above, this model is indeed more expressive than the block MDP. However, the next result answers the above question in the negative. The latent variable representation is *exponentially* weaker than the general low rank representation in the following sense:

**Proposition 2.** *For any even $n \in \mathbb{N}$, there exists an MDP that can be cast as a low rank MDP with embedding dimension $O(n^2)$, but which has $d_{\mathrm{LV}} \geq 2^{\Omega(n)}$.*

See the center and right panels of Figure 2. The result is proved by recalling that the latent variable dimension determines the non-negative rank of $T$, which can be much larger than its rank (Rothvoß, 2017; Yannakakis, 1991). It showcases how low rank MDPs are quite different from latent variable models of comparable dimension and demonstrates how embedding functions with negative values can provide significant expressiveness.

**Bellman and Witness rank.** As our last concrete connection, we remark here that the low rank MDP with a function class derived from $\Phi$ (and $\Upsilon$) admits low Bellman (resp., Witness) rank.

**Proposition 3** (Informal). *The low rank MDP model always has Bellman rank at most d. Additionally, given $\Phi$ and assuming $\phi^\star_{0:H-1} \in \Phi$, we can construct a function classes $(\mathcal{G}, \Pi)$, so that* OLIVE *when run with $(\mathcal{G}, \Pi)$ has sample complexity $\tilde{O}\left(\mathrm{poly}(d, H, K, \log|\Phi|, \epsilon^{-1})\right)$.*

See Proposition 6 for a more precise statement. An analogous result hold for the Witness rank notion of Sun et al. (2019) (see Proposition 7 in the appendix). Unfortunately both OLIVE, and the algorithm of Sun et al. (2019) are not computationally tractable, as they involve enumeration of the employed function class. We turn to the development of computationally tractable algorithms in the next section.

## 5   Main results

We now turn to the design of algorithms for representation learning and exploration in low rank MDPs. As a computational abstraction, we consider the following optimization and sampling oracles.

**Definition 3** (Computational oracles). *Define the following oracles for the classes $\Phi, \Upsilon$:*

*1. The* maximum likelihood oracle, MLE, *takes a dataset $D$ of $(x, a, x')$ triples, and returns*

$$\mathrm{MLE}(D) := \mathrm{argmax}_{\phi \in \Phi, \mu \in \Upsilon} \sum_{(x,a,x') \in D} \log(\langle \phi(x, a), \mu(x') \rangle).$$

*2. The* sampling oracle, SAMP, *is a subroutine which, for any $(\phi, \mu) \in \Phi \times \Upsilon$ and any $(x, a)$, returns a sample $x' \sim \langle \phi(x, a), \mu(\cdot) \rangle$. Multiple calls to the procedure result in independent samples.*

We assume access to both oracles as a means towards practical algorithms that avoid explicitly enumerating over all functions in $\Phi$ and $\Upsilon$. Note that related assumptions are quite common in the literature (Misra et al., 2019; Du et al., 2019b; Agarwal et al., 2014), and in practice, both oracles can be reasonably approximated whenever optimizing over $\Phi, \Upsilon$ is feasible (e.g., neural networks). Regarding MLE, other optimization oracles are possible, and in the appendix (Remark 22) we sketch how our proof can accommodate a generative adversarial oracle as a replacement (Goodfellow et al., 2014; Arora et al., 2017). While the sampling oracle is less standard, one might implement SAMP via optimization methods like the Langevin dynamics (Welling

and Teh, 2011) or through reparametrization techniques such as the Gumbel-softmax trick (Jang et al., 2017; Figurnov et al., 2018).[4] In addition, the sampling oracle can be avoided at the cost of additional real world experience, an approach we describe formally in Theorem 4 below.

## 5.1 Algorithm description

The algorithm is called FLAMBE, for "Feature Learning And Model-Based Exploration." Pseudocode is displayed in Algorithm 1. FLAMBE is iterative in nature, where in iteration $j$, we use an exploratory policy $\rho_{j-1}$ to collect a dataset of transitions, and then we pass all previously collected transitions to the MLE oracle. The optimization oracle returns embedding functions $(\hat{\phi}_h, \hat{\mu}_h)$ for each $h$ which define transition operators $\hat{T}_h$ for our updated learned model $\widehat{\mathcal{M}}$. Then FLAMBE calls a planning sub-routine to compute the exploratory policy $\rho_j$ that we use in the next iteration. After $J_{\max}$ iterations, we simply output our current model $\widehat{\mathcal{M}}$.

For the planning step, intuitively we seek an exploratory policy $\rho$ that induces good coverage over the state space when executed in the model. We do this by solving one planning problem per time step $h$ in Algorithm 2 using a technique inspired by elliptical potential arguments from linear bandits (Dani et al., 2008). Using the $h$-step model $\hat{T}_{0:h-1}$, we iteratively maximize certain quadratic forms of our learned features $\hat{\phi}_{h-1}$ to find new directions not covered by the previously discovered policies, and we update the exploratory policy to include the maximizer. The planning algorithm terminates when no policy can achieve large quadratic form, which implies that we have found all reachable directions in $\hat{\phi}_{h-1}$. This yields a *mixture policy* $\rho_h^{\mathrm{pre}}$ that is executed by sampling one of the mixture components and executing that policy for the entire episode. The component policies are linear in the learned features $\hat{\phi}_{0:h-1}$. The challenge in our analysis is to relate this coverage in the model to that in the true environment as we discuss in the next section.

Algorithm 2 is a model-based planner, so it requires no interaction with the environment. The main computational step is the optimization problem (2), which can be solved efficiently with access to the sampling oracle, essentially by running the algorithm of Jin et al. (2019) (See Lemma 15 in the appendix). Note that we are optimizing over *all policies*, which is possible because the Bellman backups in a low rank MDP are linear functions of the features (c.f., Lemma 1). The sampling oracle can also be used to approximate all expectations, and, with sufficient accuracy, this has no bearing on the final results. Our proofs do account for the sampling errors.

## 5.2 Theoretical Results

We now state the main guarantee.

**Theorem 2.** *Fix $\delta \in (0,1)$. If $\mathcal{M}$ is a low rank MDP with dimension $d$ and horizon $H$ and Assumption 1 holds, then* FLAMBE *with subroutine Algorithm 2 and appropriate settings[5] of $\beta$, $J_{\max}$, and $n$, computes a model $\widehat{\mathcal{M}}$ such that (1) holds with probability at least $1 - \delta$. The total number of trajectories collected is*

$$\tilde{O}\left(\frac{H^{22}K^9 d^7 \log(|\Phi||\Upsilon|/\delta)}{\varepsilon^{10}}\right),$$

*and the algorithm runs in polynomial time with polynomially many calls to* MLE *and* SAMP *(Definition 3).*

Thus, FLAMBE provably learns low rank MDP models in a statistically and computationally efficient manner, under Assumption 1 and Assumption 2. While the result is comparable to prior work in the dependencies on $d$, $H$, $K$ and $\varepsilon$, we instead highlight the more conceptual advances over prior work.

---

[4]We do not explicitly consider approximate oracles, but additive approximations can be accommodated in our proof. In particular, if SAMP returns a sample from a distribution that is $\varepsilon_{\mathrm{samp}}$ close in total variation to the target distribution in $\mathrm{poly}(1/\varepsilon_{\mathrm{samp}})$ time, then we retain computational efficiency.

[5]The precise settings for $\beta$, $J_{\max}$, and $n$ are given in the appendix.

---
**Algorithm 1** FLAMBE: Feature Learning And Model-Based Exploration
---
**Input:** Environment $\mathcal{M}$, function classes $\Phi, \Upsilon$, subroutines MLE and SAMP, parameters $\beta, n$.
Set $\rho_0$ to be the random policy, which takes all actions uniformly at random.
Set $D_h = \emptyset$ for each $h \in \{0, \ldots, H-1\}$.
**for** $j = 1, \ldots, J_{\max}$ **do**
    **for** $h = 0, \ldots, H-1$ **do**
        Collect $n$ samples $(x_h, a_h, x_{h+1})$ by rolling into $x_h$ with $\rho_{j-1}$ and taking $a_h \sim \text{unif}(\mathcal{A})$.
        Add these samples to $D_h$.
        Solve maximum likelihood problem: $(\hat{\phi}_h, \hat{\mu}_h) \leftarrow \text{MLE}(D_h)$.
        Set $\hat{T}_h(x_{h+1} \mid x_h, a_h) = \left\langle \hat{\phi}_h(x_h, a_h), \hat{\mu}_h(x_{h+1}) \right\rangle$.
    **end for**
    For each $h$, call planner (Algorithm 2) with $h$ step model $\hat{T}_{0:h-1}$ and $\beta$ to obtain $\rho_h^{\text{pre}}$.
    Set $\rho_j = \text{unif}(\{\rho_h^{\text{pre}} \circ \text{random}\}_{h=0}^{H-1})$, to be uniform over the discovered $h$-step policies, augmented with random actions.
**end for**
---

---
**Algorithm 2** Elliptical planner
---
**Input:** MDP $\widetilde{\mathcal{M}} = (\phi_{0:\tilde{h}}, \mu_{0:\tilde{h}})$, subroutine SAMP, parameter $\beta > 0$. Initialize $\Sigma_0 = I_{d \times d}$.
**for** $t = 1, 2, \ldots,$ **do**
    Compute (see text for details)

$$\pi_t = \underset{\pi}{\operatorname{argmax}} \, \mathbb{E}\left[\phi_{\tilde{h}}(x_{\tilde{h}}, a_{\tilde{h}})^\top \Sigma_{t-1}^{-1} \phi_{\tilde{h}}(x_{\tilde{h}}, a_{\tilde{h}}) \mid \pi, \widetilde{\mathcal{M}}\right]. \qquad (2)$$

    If the objective is at most $\beta$, halt and output $\rho = \text{unif}(\{\pi_\tau\}_{\tau < t})$.
    Compute $\Sigma_{\pi_t} = \mathbb{E}\left[\phi_{\tilde{h}}(x_{\tilde{h}}, a_{\tilde{h}})\phi_{\tilde{h}}(x_{\tilde{h}}, a_{\tilde{h}})^\top \mid \pi, \widetilde{\mathcal{M}}\right]$. Update $\Sigma_t \leftarrow \Sigma_{t-1} + \Sigma_{\pi_t}$.
**end for**
---

- The key advancement over the block MDP algorithms (Du et al., 2019b; Misra et al., 2019) is that FLAMBE applies to a significantly richer class of models with comparable function approximation assumptions. A secondary, but important, improvement is that FLAMBE does not require any reachability assumptions, unlike these previous results. We remark that Feng et al. (2020) avoid reachability restrictions in block MDPs, but their function approximation/oracle assumptions are much stronger than ours.

- Over Jin et al. (2019); Yang and Wang (2019b), the key advancement is that we address the representation learning setting where the embeddings $\phi_{0:H-1}^\star$ are not known a priori. On the other hand, our bound scales with the number of actions $K$. We believe that additional structural assumptions on $\Phi$ are required to avoid the dependence on $K$ in the representation learning setting.

- Over Jiang et al. (2017); Sun et al. (2019), the key advancement is computational efficiency. However, the low rank MDP is less general than what is covered by their theory, and our sample complexity is worse in the polynomial factors.

As remarked earlier, the logarithmic dependence on the sizes of $\Phi, \Upsilon$ can be relaxed to alternative notions of capacity for continuous classes.

We also state a sharper bound for a version of FLAMBE that operates directly on the simplex factorization. The main difference is that we use a conceptually simpler planner (See Algorithm 3 in the appendix) and the sample complexity bound scales with $d_{\text{LV}}$.

**Theorem 3.** *Fix $\delta \in (0, 1)$. If $\mathcal{M}$ admits a simplex factorization with embedding dimension $d_{\mathrm{LV}}$, Assumption 1 holds, and all $\phi \in \Phi$ satisfy $\phi(x, a) \in \Delta([d_{\mathrm{LV}}])$, then* FLAMBE *with Algorithm 3 as the subroutine and appropriate setting[6] of $J_{\max}$ and $n$ computes a model $\widehat{\mathcal{M}}$ such that (1) holds with probability at least $1 - \delta$. The total number of trajectories collected is*

$$\tilde{O}\left(\frac{H^{11}K^5 d_{\mathrm{LV}}^5 \log(|\Phi||\Upsilon|/\delta)}{\varepsilon^3}\right).$$

*The algorithm runs in polynomial time with polynomially many calls to* MLE *and* SAMP *(Definition 3).*

This bound scales much more favorably with $H, K$ and $\varepsilon$, but incurs a polynomial dependence on the latent variable dimension $d_{\mathrm{LV}}$, instead of the embedding dimension $d$. For many problems, including block MDPs, we expect that $d_{\mathrm{LV}} \approx d$, in which case using this version of FLAMBE may be preferable. However, Theorem 3 requires that we encode simplex constraints into our function class $\Phi$, for example using the softmax. When $d_{\mathrm{LV}}$ is small, this may be a practically useful design choice.

**Planning in the real world.** As our final result, we consider replacing the model-based planning subroutine with one that collects trajectories from the environment. This allows us to avoid using the sampling oracle, but our analysis requires one additional assumption. Recall the latent variable representation of Definition 2 and the fact that we can augment the trajectories with the latent variables, and let $\mathcal{Z}_h$ denote the latent state space for $T_h$, i.e., the values that $z_{h+1}$ can take. Our *reachability* assumption posits that for the MDP $\mathcal{M}$, the latent variables can be reached with non-trivial probability.

**Assumption 2** (Reachability). *There exists a constant $\eta_{\min} > 0$ such that*

$$\forall h \in \{0, \dots, H-1\}, z \in \mathcal{Z}_h : \quad \max_{\pi} \mathbb{P}\left[z_{h+1} = z \mid \pi, \mathcal{M}\right] \geq \eta_{\min}.$$

The assumption generalizes prior reachability assumptions in block MDPs (Du et al., 2019b; Misra et al., 2019), where the latent variables are referred to as "latent states." Note that reachability does not eliminate the exploration problem, as a random walk may still visit a latent variable with exponentially small probability. However, reachability is a limitation that unfortunately imposes an upper bound on the latent variable dimension $d_{\mathrm{LV}}$, as formalized in the next proposition.

**Proposition 4.** *If MDP $\mathcal{M}$ has rank $d$ and satisfies Assumption 2, then for each $h$, the latent variable dimension of $T_h$ satisfies $d_{\mathrm{LV}} \leq {dK^2}/{\eta_{\min}^2}$.*

With this new assumption, we have the following theorem, which is proved in Appendix D.

**Theorem 4.** *Fix $\delta \in (0, 1)$. In the setup of Theorem 3 with Assumption 2, a version of* FLAMBE *(described in Appendix D) computes a model $\widehat{\mathcal{M}}$ such that (1) holds with probability $1 - \delta$. The algorithm collects* $\mathrm{poly}(d_{\mathrm{LV}}, H, K, 1/\eta_{\min}, 1/\varepsilon, \log(|\Phi||\Upsilon|/\delta))$ *trajectories and runs in polynomial time with $H$ calls to* MLE.

Importantly, this instantiation of FLAMBE does not require that the function classes support efficient sampling, i.e., we do not use SAMP. On the other hand, the sample complexity degrades in comparison with our results using model-based planners. We also note that the result considers simplex representations as in Theorem 3 because (a) the calculations are considerably simpler, and (b) in light of Proposition 4 handling general representations under Assumption 2 only incurs a polynomial overhead. We believe that extending the result to accommodate general representations directly is possible.

---

[6]This version does not require the parameter $\beta$.

**Challenges in the analysis.** We highlight three main challenges in the analysis. The first challenge arises in the analysis of the model learning step, where we want to show that our model $\hat{T}_h$, learned by maximum likelihood estimation, accurately approximates the true dynamics $T_h$ on the data distributions induced by the exploratory policies $\rho_{0:j-1}$. This requires a generalization argument, but both the empirical MLE objective and population version — the KL divergence — are unbounded, so we cannot use standard uniform convergence techniques. Instead, we employ (and slightly adapt) results from the statistics literature (Van de Geer, 2000; Zhang, 2006) to show that MLE yields convergence in the Hellinger and total variation distances. While these arguments are well-known in the statistics community, we highlight them here because we believe they may be broadly useful in the context of model-based RL.

The second challenge is to transfer the model error guarantee from the exploratory distributions to distributions induced by other policies. Intuitively, error transfer should be possible if the exploratory policies cover the state space, but we must determine how we measure and track coverage. Leveraging the low rank dynamics, we show that if a policy $\pi$ induces a distribution over true features $\phi_{h-1}^\star$ that is in the span of the directions visited by the exploratory policies, then we can transfer the MLE error guarantee for $\hat{T}_h$ to $\pi$'s distribution. This suggests measuring coverage for time $h$ in terms of the second moment matrix of the true features at the previous time induced by the exploratory policies $\rho$, that is $\Sigma_{h-1} = \mathbb{E}_\rho \phi_{h-1}^\star \phi_{h-1}^{\star,\top}$. Using these matrices, we prove a sharp *simulation lemma* that bounds the model error observed by a policy $\pi$ in terms of the model error on the exploratory distribution and the probability that $\pi$ visits features for which $\phi_{h-1}^{\star,\top} \Sigma_{h-1}^{-1} \phi_{h-1}^\star$ is large.

The simulation lemma suggests a planning strategy and an "explore-or-terminate" argument: the next exploratory policy should maximize $\phi_{h-1}^{\star,\top} \Sigma_{h-1}^{-1} \phi_{h-1}^\star$, in which case either we visit some new direction and make progress, or we certify (1), since no policy can make this quantity large. However, we cannot maximize this objective directly, since we do not know $\phi_{h-1}^\star$ or $\Sigma_{h-1}$! Moreover, we cannot use $\hat{\phi}_{h-1}$ to approximate $\phi_{h-1}^\star$ in the objective, since even if the model is accurate, the features may not be. Instead we plan to find a policy that induces a well-conditioned covariance matrix in terms of $\hat{\phi}_{h-2}$, the learned features at time $h-2$. By composing this policy with a random action and applying our simulation lemma, we can show that this policy either explores at some previous time or approximately maximizes $\phi_{h-1}^{\star,\top} \Sigma_{h-1}^{-1} \phi_{h-1}^\star$, which allows us to apply the explore-or-terminate argument. Combining this reasoning with an elliptical potential argument, we can bound the number of iterations in which exploration can happen, which leads to the final result.

## 6  Discussion

This paper studies representation learning and exploration for low rank MDPs. We provide an intuitive interpretation of these models in terms of a latent variable representation, and we prove a number of structural results certifying that low rank MDPs are significantly more expressive than models studied in prior work. We also develop FLAMBE, a computationally and statistically efficient model based algorithm for system identification in low rank MDPs. Policy optimization follows as a corollary.

Our results raise a number of promising directions for future work. On the theoretical side, can we develop provably efficient model-free algorithms for representation learning in the low rank MDP? On the empirical side, can we leverage the algorithmic insights of FLAMBE to develop practically effective representation learning algorithms for complex reinforcement learning tasks? We look forward to answering these questions in future work.

# A Proofs for the structural results

In this appendix we provide proofs for the structural results in the paper. We first provide the proof of Lemma 1. In Appendix A.2 we focus on results relating to realizability and reachability. Then in Appendix A.3 we turn to the separation results of Proposition 1 and Proposition 2. Finally in Appendix A.4 we provide details about the connection to the Bellman and Witness rank.

## A.1 Proof of Lemma 1

*Proof of Lemma 1.* Fix $h$ and $V : \mathcal{X} \to [0,1]$. We drop the dependence on $h$ from the notation, with $x, a$ always corresponding to states and actions at time $h$ and $x'$ corresponding to an action at time $h+1$. Observe that as $\widehat{\mathcal{M}}$ is a low rank MDP, we have

$$\forall x, a : \quad \mathbb{E}\left[V(x') \mid x, a, \widehat{\mathcal{M}}\right] = \left\langle \hat{\phi}(x,a), \int \hat{\mu}(x')V(x') \right\rangle =: \left\langle \hat{\phi}(x,a), \theta \right\rangle$$

Combining with (1), we have that for any policy $\pi$:

$$\mathbb{E}\left[\left|\left\langle \hat{\phi}(x,a), \theta\right\rangle - \mathbb{E}\left[V(x') \mid x, a, \mathcal{M}\right]\right| \mid \pi, \mathcal{M}\right]$$
$$= \mathbb{E}\left[\left|\mathbb{E}\left[V(x') \mid x, a, \widehat{\mathcal{M}}\right] - \mathbb{E}\left[V(x') \mid x, a, \mathcal{M}\right]\right| \mid \pi, \mathcal{M}\right]$$
$$\leq \mathbb{E}\left[\left\|\left\langle \hat{\phi}(x,a), \hat{\mu}(\cdot)\right\rangle - T(\cdot \mid x, a)\right\|_{\mathrm{TV}} \mid \pi, \mathcal{M}\right] \leq \varepsilon. \qquad \square$$

## A.2 On realizability and reachability

**Proposition 5.** *Fix $M \in \mathbb{N}$, $n \leq M/2$, and any algorithm. There exists a low rank MDP over $M$ states with rank 2 and horizon 2 such that, if the algorithm collects $n$ trajectories and outputs a policy $\hat{\pi}$, then with probability at least $1/8$, $\hat{\pi}$ is at least $1/8$-suboptimal for the MDP.*

The result shows that if the low rank MDP has $M$ states, then we require $n = \Omega(M)$ samples to find a near-optimal policy with moderate probability. Thus low rank structure alone is not sufficient to obtain sample complexity guarantees that are independent of the number of states.

*Proof of Proposition 5.* The result is obtained by embedding a binary classification problem into a low rank MDP and appealing to a standard binary classification lower bound argument. We construct a family of one-step transition operators, all of which have rank 2. The state space at the current time is of size $M$ and there are two actions $\mathcal{A} := \{0, 1\}$. From each $(x, a)$ pair we transition deterministically either to $x_g$ or $x_b$, and we receive reward 1 from $x_g$ and reward 0 from $x_b$.

Formally, we denote the states as $\{x_1, \ldots, x_M\}$ and index each instance by a binary vector $v \in \{0, 1\}^M$, which specifies the good action for each state. The transition operator is

$$T_v(\cdot \mid x_j, a) = \begin{cases} x_g \text{ if } a = v_j \\ x_b \text{ if } a \neq v_j \end{cases}$$

There are therefore $2^M$ instances. Note that as there are only two states at the next time, we trivially see that that transition operator of each instance is rank 2 and the linear MDP representation is:

$$\phi_v^\star(x_j, a) = (\mathbf{1}\{a = v_j\}, \mathbf{1}\{a \neq v_j\}) \quad \text{and} \quad \mu_v^\star(x') = (\mathbf{1}\{x' = x_g\}, \mathbf{1}\{x' = x_b\}).$$

The starting distribution is uniform over $[M]$, so that in $n$ episodes, the agent collects a dataset $\{(x^{(i)}, a^{(i)}, y^{(i)}\}_{i=1}^n$ where $x^{(i)} \sim \mathrm{unif}(x_1, \ldots, x_M)$, $a^{(i)}$ is chosen by the agent and $y^{(i)}$ denotes whether

the agent transitions to $x_g$ or $x_b$. Information theoretically, this is equivalent to obtaining $n$ samples from the following data generating process: sample $j \in \mathrm{unif}([M])$ and reveal $v_j$.

In this latter process, we can apply a standard binary classification lower bound argument. Let $P_v$ denote the data distribution where indices $j$ are sampled uniformly at random and labeled by $v_j$. Let $P_v^{(n)}$ denote the product measure where $n$ samples are generated iid from $P_v$. By randomizing the instance, for any example that does not appear in the sample, the probability of error is $1/2$. Therefore the probability of error for any classifier is

$$\max_v \mathbb{E}_{S \sim P_v^{(n)}} \mathbb{P}_j[\hat{f}(j) \neq v_j] \geq \mathbb{E}_{v \sim \mathrm{Unif}(\{0,1\}^M)} \mathbb{E}_{S \sim P_v^{(n)}} \mathbb{P}_j[\hat{f}(j) \neq v_j]$$

$$= \frac{1}{M} \sum_{j=1}^M \mathbb{E}_v \mathbb{E}_{S \sim P_v^{(n)}} \mathbf{1}\{\hat{f}(j) \neq v_j\} \geq \frac{1}{M} \sum_{j=1}^M \frac{1}{2} \mathbb{P}_S[j \notin S]$$

$$= \frac{1}{2} \left(1 - \frac{1}{M}\right)^n \geq \frac{1}{2} \left(1 - n/M\right).$$

The second inequality uses the fact that if $j$ does not appear in the sample then $v_j \sim \mathrm{Ber}(1/2)$. Equivalently, we can first sample $n$ unlabeled indices, then commit to the label just on these indices, so that the label for any index not in the sample remains random. Thus for any classifier, there exists some instance for which on average over the sample, the probability of error is at least $1/4$ as long as $n \leq M/2$. This also implies that with constant probability over the sample the error rate is at least $1/8$, since for any random variable $Z$ taking values in $[0, 1]$, we have

$$\mathbb{E}[Z] \leq 1/8 \left(1 - \mathbb{P}[Z > 1/8]\right) + \mathbb{P}[Z \geq 1/8] \leq 1/8 + \mathbb{P}[Z \geq 1/8].$$

Taking $Z = \mathbb{P}_j[\hat{f}(j) \neq v_j]$, we have

$$\mathbb{P}_{S \sim P_v^{(n)}} \left[\mathbb{P}_j[\hat{f}_j \neq v_j] \geq 1/8\right] \geq \frac{1}{2} \left(1 - n/M\right) - 1/8 \geq 1/8,$$

where the last inequality holds with $n \leq M/2$.

Now, notice that we can identify any predictor with a policy in the obvious way and also that the suboptimality for a policy is precisely the classification error for the predictor. With this correspondence, we obtain the result. $\qquad\square$

*Proof of Proposition 4.* Fix stage $h$. Assume that $X := |\mathcal{X}|, d_{\mathrm{LV}} := |\mathcal{Z}_h|$ are finite, where $\mathcal{Z}_h$ is the latent state space associated with $T_h$. Recall that $\psi_h(x, a) \in \Delta(\mathcal{Z})$ maps each state action pair to a distribution over latent states. As $X, |\mathcal{A}|, d_{\mathrm{LV}}$ are all finite, we may collect these vectors as a matrix $\Psi \in \mathbb{R}^{XK \times d_{\mathrm{LV}}}$ with rows corresponding to $\psi_h(x, a)$. A policy $\pi$ induces a distribution over $(x, a)$ pairs, which we call $p_{\pi,h} \in \Delta(\mathcal{X} \times \mathcal{A})$. The corresponding distribution over latent variables at stage $h$ is therefore $p_{\pi,h}^\top \Psi$.

We re-express $p_{\pi,h}$ in two steps. First, we can write $p_{\pi,h} = A_{\pi,h} \times \tilde{p}_{\pi,h}$ where $A_{\pi,h} \in \mathbb{R}^{XK \times X}$ is a matrix where the $x^{\mathrm{th}}$ column describes the distribution $\pi(\cdot \mid x) \in \Delta(\mathcal{A})$, and $\tilde{p}_{\pi,h} \in \Delta(\mathcal{X})$ is the distribution over $x_h$ induced by policy $\pi$. Note that $A_{\pi,h}$ is column stochastic (it is non-negative with each column summing to 1). In fact it has additional structure, since in column $x$, only the rows corresponding to $(x, \cdot)$ are non-zero, but this will not be essential for our arguments. As $A_{\pi,h}^\top$ is therefore row-stochastic and the product of two row-stochastic matrices is also row-stochastic, we have that $A_{\pi,h}^\top \Psi \in \mathbb{R}^{X \times d_{\mathrm{LV}}}$ is also row-stochastic.

Next, we use the dynamics at stage $h - 1$ to re-write $\tilde{p}_{\pi,h}$, which is the state distribution induced by policy $\pi$ at stage $h$. As $T_{h-1}$ is also rank $d$, we can write $T(x_h \mid x_{h-1}, a_{h-1}) = \langle \phi_{h-1}(x_{h-1}, a_{h-1}), \mu_{h-1}(x_h) \rangle$, and we can collect the embeddings $\mu_{h-1}(x_h)$ as columns of a $d \times X$ matrix $U_{h-1}$. With these definitions,

$$\tilde{p}_{\pi,h} = \mathbb{E}\left[U_{h-1}^\top \phi_{h-1}(x_{h-1}, a_{h-1}) \mid \pi, \mathcal{M}\right] = U_{h-1}^\top v_{\pi,h-1}.$$

12

Here $\mathcal{M}$ is the MDP in consideration. In summary, for any policy $\pi$, the distribution over latent variable $z_{h+1}$ (which generates $x_{h+1}$) induced by policy $\pi$ can be written as

$$\mathbb{P}\left[z_{h+1} = \cdot \mid \pi, \mathcal{M}\right] = v_{\pi,h-1}^\top U_{h-1} A_{\pi,h}^\top \Psi \in \mathbb{R}^{d_{\mathrm{LV}}}.$$

Now, let us use our linear-algebraic re-writing to express the reachability condition. If a latent variable $z \in \mathcal{Z}_h$ is reachable, then there exists some policy $\pi_z$ such that $\mathbb{P}[z_{h+1} = z \mid \pi_z, \mathcal{M}] \geq \eta_{\min}$. First of all, by importance weighting on the last action $a_h$, we have:

$$\mathbb{P}\left[z_{h+1} = z \mid a_{0:h-1} \sim \pi_z, a_h \sim \mathrm{unif}(\mathcal{A}), \mathcal{M}\right] \geq \eta_{\min}/K.$$

The normalization condition on $\phi_{h-1}$ leads to the upper bound

$$\frac{\eta_{\min}}{K} \leq \left| v_{\pi_z,h-1}^\top U_{h-1} A_h^\top \Psi e_z \right| \leq \max_{v:\|v\|_2 \leq 1} \left| v^\top U_{h-1} A_h \Psi e_z \right| = \|U_{h-1} A_h \Psi e_z\|_2 .$$

where we use $A_h$ to denote the action-selection matrix induced by the uniform policy.

Next, consider some $\ell_\infty$-bounded vector $w \in \mathbb{R}^{\mathcal{Z}}$ with $\|w\|_\infty \leq 1$. The fact that $A_h^\top \Psi$ is row-stochastic implies that

$$\left\| A_h^\top \Psi w \right\|_\infty \leq \max_{x,a} \left| \psi(x,a)^\top w \right| \leq \|w\|_\infty .$$

Therefore, using the normalization condition on $U_{h-1}$ we have

$$\max_{w:\|w\|_\infty \leq 1} \left\| U_{h-1} A_h^\top \Psi w \right\|_2^2 \leq \max_{w:\|w\|_\infty \leq 1} \|U_{h-1} w\|_2^2 \leq d.$$

We will select a vector $w \in \{\pm 1\}^{d_{\mathrm{LV}}}$, for which we know this upper bound holds. We select the vector iteratively, peeling off latent variables that are reachable. For brevity, define $B := U_{h-1} A_h^\top \Psi$ and observe that

$$\|Bw\|_2^2 = \|Be_{z_1} w[z_1]\|_2^2 + 2\langle Be_{z_1} w[z_1], \sum_{z \neq z_1} Be_z w[z] \rangle + \| \sum_{z \neq z_1} Be_z w[z]\|_2^2.$$

If $z_1$ is $\eta_{\min}$ reachable, then the first term is at least $(\eta_{\min}/K)^2$ by the above calculation and the fact that we take $w[z_1] \in \{\pm 1\}$. Then we ensure that the cross-term is non-negative by setting $w[z_1]$ appropriately. Note that $w[z_1]$ is formally a function of the remaining coordinates of $w$, but we have not introduced any constraint on these remaining coordinates. Therefore, for $z_1$ is $\eta_{\min}$ reachable, we get (for this partially specified $w$)

$$\|Bw\|_2^2 \geq (\eta_{\min}/K)^2 + \| \sum_{z \neq z_1} Be_z w[z]\|_2^2$$

Continuing in this way, we iteratively peel off latent variables that are $\eta_{\min}$ reachable, and for each we gain $(\eta_{\min}/K)^2$ in the lower bound. Therefore, if all $d_{\mathrm{LV}}$ latent variables are reachable, there exists some $w \in \{\pm 1\}^{d_{\mathrm{LV}}}$ such that

$$d_{\mathrm{LV}} \cdot (\eta_{\min}/K)^2 \leq \left\| U_{h-1}^\top A_h^\top \Psi w \right\|_2^2 \leq d,$$

which implies that we must have $d_{\mathrm{LV}} \leq dK^2/\eta_{\min}^2$. $\qquad \square$

## A.3 Separation results

*Proof of Proposition 1.* Fix $N$ and consider a MDP with horizon 2, where at stage 1 there is only one state $x$ and two actions $a_1, a_2$. At stage 2 there are $N$ possible states, so that $T(\cdot \mid x, a_i) \in \Delta([N])$ for each $i \in \{1, 2\}$. We define the transition operator for stage 1, called $T$ for brevity, explicitly in terms of its factorization. Let $\phi(x, a_1) = e_1$ and $\phi(x, a_2) = e_2$ where $e_1, e_2 \in \mathbb{R}^2$ denotes the two standard basis elements in two dimensions. We define $\mu_1(i) = 1/N$, $\mu_2(i) = i/(\sum_{j=1}^N j)$ and $\mu(i) = (\mu_1(i), \mu_2(i)) \in \mathbb{R}^2$. Thus $T(x' = i \mid x, a) = \langle \phi(x, a), \mu(i) \rangle$, which can be easily verified to be a valid transition operator. By construction $T$ has rank 2.

For clarity we express $T$ as the $2 \times N$ matrix.

$$T := \begin{pmatrix} 1/N & 1/N & \cdots & 1/N \\ 1/(\sum_{j=1}^N j) & 2/(\sum_{j=1}^N j) & \cdots & N/(\sum_{j=1}^N j) \end{pmatrix}.$$

We now show that the block MDP representation must have $N$ latent states. Suppose the block MDP representation is $T(x' = i \mid x, a) = \langle \phi_B(x, a), \mu_B(i) \rangle$. The block MDP representation requires that for each index $i$ the vector $\mu_B(i)$ is one-sparse. From this, we deduce a constraint that arises when two states belong to the same block. If $i, j$ belong to the same block, say block $b$, then for each $(x, a) \in \mathcal{X} \times \mathcal{A}$, we have

$$T(x' = i \mid x, a) = \phi_B(x, a)[b]\mu_B(i)[b] = \frac{\mu_B(i)[b]}{\mu_B(j)[b]} \cdot \phi_B(x, a)[b]\mu_B(j)[b]$$

$$= \frac{\mu_B(i)[b]}{\mu_B(j)[b]} \cdot T(x' = j \mid x, a)$$

In words, if states $i, j$ at stage 2 belong to the same block, then the vectors $T(x' = i \mid \cdot), T(x' = j \mid \cdot)$ must be pairwise linearly dependent.[7] Based on our construction, $T(x' = i \mid \cdot) = \mu(i)$, which is just the $i^{\text{th}}$ column of the matrix $T$. By inspection, all $N$ vectors are pairwise linearly independent, and so we can conclude that the block MDP representation must have $N$ latent states. □

*Proof of Proposition 2.* We consider a one step transition operator $T$ that we instantiate to be the slack matrix describing a certain polyhedral set. Let $n$ be even and let $K_n$ be the complete graph on $n$ vertices. To set up the notation we will work with vectors $x \in \mathbb{R}^{\binom{n}{2}}$ that associate a weight to each edge. We index the vectors as $x_{u,v}$ where $u \neq v \in [n]$ correspond to vertices.

A result of Edmonds (1965) states that the perfect matching polytope, which is the convex hull of all edge-indicator vectors corresponding to perfect matchings, can be explicitly written in terms of "odd-cut" constraints:

$$\mathcal{P}_n := \text{conv} \left\{ \mathbf{1}_M \in \mathbb{R}^{\binom{n}{2}} \mid M \text{ is a perfect matching in } K_n \right\}$$

$$= \left\{ x \in \mathbb{R}^{\binom{n}{2}} : x \succeq 0, \forall v : \sum_u x_{u,v} = 1, \forall U \subset [n], |U| \text{ odd} \sum_{v \notin U} \sum_{u \in U} x_{u,v} \geq 1 \right\}.$$

This polytope has exponentially many vertices and exponentially many constraints. Formally, there are $V := \frac{n!}{2^{n/2}(n/2)!}$ vertices, corresponding to perfect matchings in $K_n$, and the number of constraints is $C := 2^{\Omega(n)}$ corresponding to the number of odd-sized subsets of $[n]$. By adding one dimension to account for the offsets in the inequality constraints, we can enumerate the vertices $v_1, \ldots, v_V \in \mathbb{R}^{\binom{n}{2}+1}$ and the constraints $c_1, \ldots, c_C \in \mathbb{R}^{\binom{n}{2}+1}$, such that $\langle c_i, v_j \rangle \geq 0$ for all $i, j$. Then, we define the *slack matrix* for this polytope to be $Z \in \mathbb{R}_+^{C \times V}$ with entries $Z_{i,j} = \langle c_i, v_j \rangle$.

---

[7]Note that this is equivalent to the notion of backward kinematic inseparability (Misra et al., 2019).

14

This slack matrix clearly has rank $\binom{n}{2} + 1 = O(n^2)$. On the other hand, we claim that the non-negative rank is at least $2^{\Omega(n)}$. This follows from (a) the fact that $\mathcal{P}_n$ has extension complexity $2^{\Omega(n)}$ (Rothvoß, 2017), (b) the extension complexity of a polytope is exactly the non-negative rank of its slack matrix (Yannakakis, 1991; Fiorini et al., 2013).

Next, we define the transition operator $T$. We associate each $(x, a)$ pair with a constraint $c_i$ and each $x'$ with a vertex $v_j$. Then we define

$$T(x' \mid x, a) = \frac{\langle c_i, v_j \rangle}{\sum_{k=1}^{V} \langle c_i, v_k \rangle}$$

This is easily seen to be a distribution for each $(x, a)$ pair. We can represent $T$ as a $C \times V$ matrix $T = DZ$ where $D$ is a diagonal matrix (with strictly positive diagonal) and $Z$ is the slack matrix defined above.

We conclude the proof with two facts from Cohen and Rothblum (1993). First, the non-negative rank is preserved under positive diagonal rescaling, and so the non-negative rank of $T$ is also $2^{\Omega(n)}$. Second, for a row-stochastic matrix $P$, the non-negative rank is equal to the smallest number of factors we can use to write $P = RS$ where both $R$ and $S$ are row-stochastic (here factors refers to the internal dimension). It is immediate that the simplex features representation corresponds to such a row-stochastic factorization, and so we see that any simplex features representation of $T$ must have embedding dimension at least $2^{\Omega(n)}$. $\square$

## A.4 On Bellman and Witness rank

We now state the formal version of Proposition 3. We consider the value-function/policy decomposition studied by Jiang et al. (2017) where we approximate the value functions with a class $\mathcal{G} : \mathcal{X} \to [0, H]$ and the policies with a class $\Pi : \mathcal{X} \to \mathcal{A}$. Given an explicit reward function $R$ with range $[0, 1]$ and the function class $\Phi$ of candidate embeddings, we define these two classes as:

$$\Pi(\Phi) := \left\{ \pi : x_h \mapsto \underset{a \in \mathcal{A}}{\operatorname{argmax}} \langle \phi_h(x_h, a), \theta_h \rangle + R(x_h, a_h) : \theta_{0:H-1} \in B_d(H\sqrt{d}), \phi_{0:H-1} \in \Phi \right\},$$

$$\mathcal{G}(\Phi) := \left\{ g : x_h \mapsto \max_a \langle \phi_h(x_h, a), \theta_h \rangle + R(x_h, a_h) : \theta_{0:H-1} \in B_d(H\sqrt{d}), \phi_{0:H-1} \in \Phi \right\}.$$

Here $B_d(\cdot)$ is the Euclidean ball in $d$ dimensions with the specified radius. We have the following proposition:

**Proposition 6.** *The low rank MDP model with* any *function classes $\mathcal{G} \subset \mathcal{X} \to [0, B]$ and $\Pi \subset \mathcal{X} \to \Delta(\mathcal{A})$ has bellman rank at most $d$ with normalization parameter $O(B\sqrt{d})$. Additionally, for any known reward function $R$ with range $[0, 1]$ and assuming $\phi_{0:H-1}^\star \in \Phi$, the optimal policy and value function lie in $(\mathcal{G}(\Phi), \Pi(\Phi))$, and so OLIVE has sample complexity $\tilde{O}\left(\operatorname{poly}(d, H, K, \log|\Phi|, \epsilon^{-1})\right)$.*

*Proof of Proposition 6.* The result is essentially Proposition 9 in Jiang et al. (2017), who address the simplex representation case. We address the general case and also verify the realizability assumption.

Consider any explicitly specified reward function $R : \mathcal{X} \times \mathcal{A} \times \{0, \ldots, H-1\} \to [0, 1]$ and any low rank MDP with embedding functions $\phi_{0:H-1}^\star, \mu_{0:H-1}^\star$ and embedding dimension $d$. For any policy $\pi, \pi'$ and any value function $g : \mathcal{X} \to \mathbb{R}$ we define the *average Bellman error* (Jiang et al., 2017) as

$$\mathcal{E}(\pi, (g, \pi'), h) := \mathbb{E}\left[g(x_h) - R_h(x_h, a_h) - g(x_{h+1}) \mid a_{0:h-1} \sim \pi, a_h = \pi'(x_h), \mathcal{M}\right],$$

We also introduce the shorthand

$$\Delta((g, \pi'), x_h) := \mathbb{E}\left[g(x_h) - R_h(x_h, a_h) - g(x_{h+1}) \mid x_h, a_h = \pi'(x_h)\right].$$

Then, in the low rank MDP, the average Bellman error admits a factorization as follows

$$\mathcal{E}(\pi, (g, \pi'), h) = \mathbb{E}\left[\Delta((g, \pi'), x_h) \mid x_h \sim \pi\right]$$

$$= \left\langle \mathbb{E}\left[\phi^\star_{h-1}(x_{h-1}, a_{h-1}) \mid \pi\right], \int \mu^\star_{h-1}(x_h)\Delta((g, \pi'), x_h)d(x_h)\right\rangle$$

$$=: \left\langle \nu_h(\pi), \xi_h((g, \pi'))\right\rangle$$

We also have the normalization $\|\nu_h(\pi)\|_2 \leq 1$ and $\|\xi_h((g, \pi))\|_2 \leq (2B+1)\sqrt{d}$. This final calculation is based on the triangle inequality, the bounds on $g$ and $R$ and the normalization of $\mu^\star_{h-1}$. Thus for any low rank MDP and *any* (bounded) function class $\mathcal{G}, \Pi$, the Bellman rank is at most $d$ with norm parameter $O(B\sqrt{d})$.

To prove that OLIVE has low sample complexity, we need to verify that the optimal policy and optimal value function lie in $\Pi(\Phi)$ and $\mathcal{G}(\Pi)$ respectively. Then we must calculate the statistical complexity of these two classes. Observe that we can express the Bellman backup of any function $V : \mathcal{X} \to \mathbb{R}$ as a linear function in the optimal embedding $\phi^\star$:

$$(\mathcal{T}_h V)(x, a) := \mathbb{E}[R_h(x, a) + V(x') \mid x, a, h] = R_h(x, a) + \left\langle \phi^\star_h(x, a), \int \mu^\star_h(x')V(x')d(x')\right\rangle$$

$$= R_h(x, a) + \langle \phi^\star_h(x, a), w \rangle.$$

for some vector $w$. Moreover, if $V : \mathcal{X} \to [0, H]$, we know that $\|w\| \leq H\sqrt{d}$. In particular, this implies that the optimal $Q$ function is a linear function in the true embedding functions $\phi^\star_{0:H-1}$, and so realizability holds for $\mathcal{G}(\Phi), \Pi(\Phi)$. These function classes have range $B = O(H\sqrt{d})$ so the normalization parameter in the Bellman rank definition is $O(Hd)$.

Finally, we must calculate the statistical complexity of these two classes. For $\Pi(\Phi)$ the Natarajan dimension is at most $\tilde{O}\left(H(d + \log|\Phi|)\right)$, since for each $h$, we choose $\phi_h$ and a $d$-dimensional linear classifier. Analogously the pseudo-dimension of $\mathcal{G}(\Phi)$ is $\tilde{O}\left(H(d + \log|\Phi|)\right)$. Formally, we give a crude upper bound on the growth function, focusing on $\Pi(\Phi)$. Fix $h$, let $S$ be a sample of $n$ pairs $(x, a)$, and let $h_1, h_2 : S \to \{0, 1\}$ such that $h_1(x, a) \neq h_2(x, a)$ for all points in the sample. Since once we fix $\phi \in \Phi$, we have a linear class, we can vary $\theta$ to match $h_1, h_2$ on at most $(n+1)^d$ subsets $T \subset S$. Then by varying $\phi \in \Phi$ we can match $h_1, h_1$ in total on $|\Phi|(n+1)^d \leq n^{O(d+\log|\Phi|)}$ subsets. If $S$ is shattered, this means that $2^n \leq n^{O(d+\log|\Phi|)}$, which means that the Natarajan dimension is $O((d + \log|\Phi|)\log(d + \log|\Phi|))$. This calculation is for a fixed $h$, but the same argument yields the bound of $\tilde{O}(H(d + \log|\Phi|))$. Instantiating, we obtain the sample complexity bound for OLIVE. $\qquad\square$

For the model-based version using the witness rank, the arguments are more straightforward.

**Proposition 7.** *The low rank MDP model with any candidate model class $\mathcal{P}$ has witness rank at most $d$, with norm parameter $O(\sqrt{d})$. Additionally, for any explicitly specified reward function $R$ with range $[0, 1]$ and under Assumption 1, the algorithm of Sun et al. (2019) (with witness class of all bounded functions) has sample complexity $\tilde{O}\left(\text{poly}(d, K, H, \log|\Phi||\Upsilon|, \varepsilon^{-1})\right)$.*

*Proof.* Given a model $M$ and an explicit reward function $R$, we use $\pi_M$ to denote the optimal policy for $R$ with transitions governed by $M$. Then, for two models $M_1, M_2$ and a time step $h$ the witness model misfit, when instantiated with the test function class as all bounded functions, is defined as

$$\mathcal{W}(M_1, M_2, h) := \mathbb{E}\left[\|M_2(\cdot \mid x_h, a_h) - \mathcal{M}(\cdot \mid x_h, a_h)\|_{\text{TV}} \mid a_{0:h-1} \sim \pi_{M_1}, a_h = \pi_{M_2}, \mathcal{M}\right].$$

Here we use the notation $M(\cdot \mid x_h, a_h)$ to denote the transition operator implied by $M$ at stage $h$. Recall that $\mathcal{M}$ is the true MDP. In words, the witness model misfit is the one-step TV error between candidate model $M_2$ and the true environment $\mathcal{M}$ on the data distribution induced by executing policy $\pi_{M_1}$ for $h$ steps.

Using the backing up argument from the proof of Proposition 6, it is easy to see that the witness model misfit admits a factorization as

$$\mathcal{W}(M_1, M_2, h) = \left\langle \mathbb{E}\left[\phi_{h-1}^\star(x_{h-1}, a_{h-1}) \mid \pi_{M_1}, \mathcal{M}\right], \int \nu_{h-1}^\star(x_h)\Delta(x_h, M_2) \right\rangle$$

where $\Delta(x_h, M_2)$ is the expected total variation distance between $M_2$ and $\mathcal{M}$ on $(x_h, \pi_{M_2}(x_h))$. Based on this calculation, the witness rank is at most $d$ and the normalization parameter is at most $O(\sqrt{d})$. It is more straightforward to see that realizability holds here, and so the algorithm of Sun et al. (2019) has the stated sample complexity. $\qquad\square$

## B   Analysis of FLAMBE

As a reminder, FLAMBE interacts with a low rank MDP $\mathcal{M}$, with time horizon $H$ and with non-stationary dynamics $T_h(x_{h+1} \mid x_h, a_h) = \langle \phi_h^\star(x_h, a_h), \mu_h^\star(x_{h+1})\rangle$. We assume that for each $h$ the operators $\phi_h^\star, \mu_h^\star$ embed into $\mathbb{R}^d$. We use the shorthand $\mathbb{E}_\pi[\cdot] = \mathbb{E}[\cdot \mid \pi, \mathcal{M}]$ to denote expectations when policy $\pi$ interacts with the real MDP $\mathcal{M}$ and $\hat{\mathbb{E}}_\pi[\cdot] = \mathbb{E}\left[\cdot \mid \pi, \widehat{\mathcal{M}}\right]$ for expectations when the policy interacts with the estimated MDP $\widehat{\mathcal{M}}$, which has dynamics $\hat{T}_{0:H-1}$. Note that this MDP model changes from iteration to iteration. When necessary we will use $\hat{\mathbb{E}}_{j,\pi}[\cdot]$ to denote the MDP model learned in the $j^{\text{th}}$ iteration of FLAMBE.

The analysis of FLAMBE is based on a potential function argument. The key quantities are the second moment matrices of the real features induced by the policies $\rho_0, \rho_1, \ldots$ at each time $h$. Formally, for $h \in \{0, \ldots, H-1\}$ and $j \in [J_{\max}]$ we define

$$\Sigma_{h,j} := \lambda I_{d\times d} + \sum_{i=0}^{j-1} \mathbb{E}_{\rho_i}\left[\phi_h^\star(x_h, a_h)\phi_h^\star(x_h, a_h)^\top\right],$$

where $\lambda > 0$ is a small constant we will set towards the end of the proof. Note that $\Sigma_{h,j} \succ 0$ for all $h, j$.

The importance of $\Sigma_{h,j}$ is demonstrated in the next result, which establishes an accuracy guarantee for the model $\hat{T}_{0:H-1}$ learned in iteration $j$. The result is a corollary of Theorem 21.

**Corollary 5.** *Fix $j \geq 1$, $h \in \{1, \ldots H-1\}$, $\delta \in (0, 1)$, and let $\rho_0, \ldots, \rho_{j-1}$ be any (possibly data-dependent) policies, with $\Sigma_{h,j}$ defined accordingly. Let $D_h$ be a dataset of $nj$ examples where for each $0 \leq i < j$ we collect $n$ triples $(x_h, a_h, x_{h+1})$ by rolling in with $\rho_i$ to $x_h$ and taking $a_h$ uniformly at random. Then with probability $1 - \delta$ the output $(\hat{\phi}_h, \hat{\mu}_h)$ of $\mathrm{MLE}(D_h)$ satisfies*

$$\left\| \int \mu_{h-1}(x_h)\mathrm{unif}(a_h) \left\| \left\langle \hat{\phi}_h(x_h, a_h), \hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h, a_h)\right\|_{\mathrm{TV}} \right\|_{\Sigma_{h-1,j}}^2 \leq \lambda d + \frac{2\log(|\Phi||\Upsilon|/\delta)}{n}.$$

*Additionally, for any $j \geq 1$, with probability at least $1 - \delta$ we have*

$$\left\| \left\langle \hat{\phi}_0(x_0, a_0), \hat{\mu}_0(\cdot)\right\rangle - T(\cdot \mid x_0, a_0)\right\|_{\mathrm{TV}}^2 \leq \frac{2\log(|\Phi||\Upsilon|/\delta)}{n}.$$

*Proof.* For shorthand, we use $v_h$ to denote the $d$-dimensional vector on the left hand side of the desired

bound. Then, the left hand side is

$$\|v_h\|_{\Sigma_{h-1},j}^2 = \lambda \|v_h\|_2^2 + \sum_{i=0}^{j-1} \mathbb{E}_{\rho_i}\left[\left(\phi_{h-1}^\star(x_{h-1},a_{h-1})^\top v_h\right)^2\right]$$

$$= \lambda \|v_h\|_2^2 + \sum_{i=0}^{j-1} \mathbb{E}_{\rho_i}\left[\left(\mathbb{E}\left[\left\|\left\langle\hat{\phi}_h(x_h,a_h),\hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h,a_h)\right\|_{\mathrm{TV}} \mid x_{h-1},a_{h-1}\right]\right)^2\right]$$

$$\leq \lambda d + \sum_{i=0}^{j-1} \mathbb{E}_{\rho_i}\left[\left\|\left\langle\hat{\phi}_h(x_h,a_h),\hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h,a_h)\right\|_{\mathrm{TV}}^2 \mid a_h \sim \mathrm{unif}(\mathcal{A})\right].$$

The first term appears in the desired bound, so now we focus on the second term. We have $nj$ total examples that form a martingale process, since $\rho_i$ depends on all of the data collected in previous iterations. Applying Theorem 21, we see that with probability $1 - \delta$:

$$\sum_{i=0}^{j-1} n \cdot \mathbb{E}_{\rho_i}\left[\left\|\left\langle\hat{\phi}_h(x_h,a_h),\hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h,a_h)\right\|_{\mathrm{TV}}^2 \mid a_h \sim \mathrm{unif}(\mathcal{A})\right] \leq 2\log(|\Phi||\Upsilon/\delta|),$$

where the factor of $n$ arises since we collect $n$ examples from $\rho_i$. Re-arranging we obtain the first bound. The bound for $h = 0$ is a direct application of Theorem 21, since we assume there is a fixed starting state $x_0$ with a single available action. $\qquad\square$

Now that we have established an accuracy guarantee in terms of the previous exploratory policies, we state and prove the main technical "simulation" lemma. The following notation is helpful. Given an MDP model $\hat{\phi}_{0:H-1}, \hat{\mu}_{0:H-1}$ and positive definite matrices $\Sigma_0,\ldots,\Sigma_{H-1}$, define

$$\forall h \geq 1 : \mathrm{err}_h(\Sigma_{h-1}) := \left\|\int \mu_{h-1}(x_h)\mathrm{unif}(a_h)\left\|\left\langle\hat{\phi}_h(x_h,a_h),\hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h,a_h)\right\|_{\mathrm{TV}}\right\|_{\Sigma_{h-1}}^2,$$

$$\mathrm{err}_0 := \left\|\left\langle\hat{\phi}_0(x_0,a_0),\hat{\mu}_0(\cdot)\right\rangle - T_0(\cdot \mid x_0,a_0)\right\|_{\mathrm{TV}}^2.$$

Further, for each $h \geq 1$, define $\mathcal{K}_h(\Sigma_h) := \left\{(x,a) \in \mathcal{X} \times \mathcal{A} : \|\phi_h^\star(x_h,a_h)\|_{\Sigma_h^{-1}}^2 \leq 1\right\}$. Let $M_{\mathcal{K}}$ be the MDP with non-stationary transition operator $T_{h,\mathcal{K}}$ defined as

$$T_{h,\mathcal{K}}(x_{h+1} \mid x_h,a_h) = \begin{cases} \langle\phi_h^\star(x_h,a_h),\mu_h^\star(x_{h+1})\rangle & \text{if } (x_h,a_h) \in \mathcal{K}_h(\Sigma_h) \\ \mathbf{1}\{x_{h+1} = x_{\mathrm{absorb}}\} & \text{if } (x_h,a_h) \notin \mathcal{K}_h(\Sigma_h) \end{cases},$$

where $x_{\mathrm{absorb}}$ is a special self-looping absorbing state with a single action $a_{\mathrm{absorb}}$ such that $T(x_{\mathrm{absorb}} \mid x_{\mathrm{absorb}},a_{\mathrm{absorb}}) = 1$ always. The initial transition $T_{0,\mathcal{K}}$ is identical to $T_0$. The intuition is that $\mathcal{K}$ denotes the set of "known" state-action pairs, and the MDP $M_{\mathcal{K}}$ terminates any episode that escapes the known set. In all of these definitions, we suppress the dependence on $\Sigma_h$ when it is clear from context. We always consider $(x_{\mathrm{absorb}}, a_{\mathrm{absorb}})$ to be *known*.

**Lemma 6.** *Let $\hat{\phi}_{0:H-1}, \hat{\mu}_{0:H-1}$ be an MDP model and let $\Sigma_{0:H-1}$ be positive definite matrices. Assume that*

$$\forall h \in \{0,\ldots,H-1\} : \mathrm{err}_h(\Sigma_{h-1}) \leq \varepsilon_{\mathrm{TV}}.$$

*Let $f : \mathcal{X} \times \mathcal{A} \to [0,1]$ be any function such that $f(x_{\mathrm{absorb}},a_{\mathrm{absorb}}) = 0$, and let $\pi$ be any policy. Then for any $h \in \{0,\ldots,H-1\}$*

$$\mathbb{E}_\pi\left[f(x_h,a_h) \mid M_{\mathcal{K}}\right] - HK\sqrt{\varepsilon_{\mathrm{TV}}} \leq \hat{\mathbb{E}}_\pi\left[f(x_h,a_h)\right] \leq \mathbb{E}_\pi\left[f(x_h,a_h) \mid M_{\mathcal{K}}\right] + HK\sqrt{\varepsilon_{\mathrm{TV}}}$$

$$+ \sum_{h'=0}^{h-1} \mathbb{P}\left[(x_{h'},a_{h'}) \notin \mathcal{K}_{h'} \mid \pi, M_{\mathcal{K}}\right].$$

This lemma establishes a sharp relationship between the learned MDP $\widehat{\mathcal{M}}$ and an *absorbing* MDP $\mathcal{M}_\mathcal{K}$, defined in terms of the matrices $\Sigma_h$, which also governs the estimation error for $\widehat{\mathcal{M}}$. Intuitively the error guarantee implies that $\widehat{\mathcal{M}}$ closely approximates $\mathcal{M}_\mathcal{K}$ provided we stay within the known set $\mathcal{K}_{0:H-1}$. Conversely the difference in value between the two MDPs can be bounded in terms of the escaping probability, which is the third term on the right hand side of the bound.

Note also that the above lemma, with $\varepsilon_{\mathrm{TV}} = 0$, can be used to compare $\mathcal{M}$ with $\mathcal{M}_\mathcal{K}$, which yields that for any non-negative function $f$ and any policy $\pi$:

$$\mathbb{E}_\pi\left[f(x_h, a_h) \mid \mathcal{M}_\mathcal{K}\right] \le \mathbb{E}_\pi\left[f(x_h, a_h)\right] \le \mathbb{E}_\pi\left[f(x_h, a_h) \mid \mathcal{M}_\mathcal{K}\right] + \sum_{h=0}^{h-1} \mathbb{P}\left[(x_{h'}, a_{h'}) \notin \mathcal{K}_{h'} \mid \pi, \mathcal{M}_\mathcal{K}\right]. \quad (3)$$

This bound actually holds for any sets $\mathcal{K}_h$. We now turn to the proof of Lemma 6.

*Proof.* Let $\hat{V}_{h'}(x) := \hat{\mathbb{E}}_\pi[f(x_h, a_h) \mid x_{h'} = x]$ denote the value function (relative to $f$) in the model and let $V_{h',\mathcal{K}}(x)$ denote the analogous quantity in the absorbing MDP. We have the following telescoping identity:

$$\hat{\mathbb{E}}_\pi[f(x_h, a_h)] - \mathbb{E}_\pi\left[f(x_h, a_h) \mid M_\mathcal{K}\right] = \int \left(\hat{T}_0(x_1 \mid x_0, a_0) - T_{0,\mathcal{K}}(x_1 \mid x_0, a_0)\right) \hat{V}_1(x_1)$$

$$+ \mathbb{E}_\pi\left[\hat{V}_1(x_1) - V_{1,\mathcal{K}}(x_1) \mid M_\mathcal{K}\right]$$

$$= \sum_{h'=0}^{h-1} \mathbb{E}_\pi\left[\int (\hat{T}_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}) - T_{h',\mathcal{K}}(x_{h'+1} \mid x_{h'}, a_{h'}))\hat{V}_{h'+1}(x_{h'+1}) \mid M_\mathcal{K}\right].$$

Note that $\hat{V}_h(x) = V_{h,\mathcal{K}}(x)$ at the time $h$ where we apply function $f$. This means that we only accumulate errors up to time $h - 1$. We now work with one of these terms. In the remainder of the proof, unless otherwise specified, all expectations are taken by executing $\pi$ in $M_\mathcal{K}$. By adding and subtracting $T_{h'}$, we get two terms

$$\mathrm{Term1}_{h'} := \mathbb{E}\left[\int (\hat{T}_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}) - T_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}))\hat{V}_{h'+1}(x_{h'+1})\right]$$

$$\mathrm{Term2}_{h'} := \mathbb{E}\left[\int (T_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}) - T_{h',\mathcal{K}}(x_{h'+1} \mid x_{h'}, a_{h'}))\hat{V}_{h'+1}(x_{h'+1})\right].$$

For $\mathrm{Term1}_{h'}$, note that the expression evaluates to zero if $x_{h'} = x_{\mathrm{absorb}}$ since both $\hat{T}_{h'}$ and $T_{h'}$ agree that $x_{\mathrm{absorb}}$ has a single self-looping action. We now bound Term1 at time $h' = 1$, although exactly the same argument applies to $h' > 0$. Defining $\mathrm{err}(x_1, a_1) := \left\|\hat{T}_1(\cdot \mid x_1, a_1) - T_1(\cdot \mid x_1, a_1)\right\|_{\mathrm{TV}}$ and by applying Holder's inequality, we have

$$\mathrm{Term1}_1 \le \mathbb{E}\left[\mathbf{1}\{x_1 \ne x_{\mathrm{absorb}}\}\left\|\hat{T}_1(\cdot \mid x_1, a_1) - T_1(\cdot \mid x_1, a_1)\right\|_{\mathrm{TV}}\right]$$

$$\le K \cdot \mathbb{E}\left[\mathbf{1}\{x_1 \ne x_{\mathrm{absorb}}\}\mathrm{unif}(a_1)\mathrm{err}(x_1, a_1)\right]$$

$$= K \cdot \mathbb{E}\left[\phi_0^\star(x_0, a_0)\mathbf{1}\left\{\|\phi_0^\star(x_0, a_0)\|_{\Sigma_0^{-1}}^2 \le 1\right\}\right] \cdot \int \mu_0^\star(x_1)\mathrm{unif}(a_1)\mathrm{err}(x_1, a_1)$$

$$\le K \cdot \mathbb{E}\left[\|\phi_0^\star(x_0, a_0)\|_{\Sigma_0^{-1}}\mathbf{1}\left\{\|\phi_0^\star(x_0, a_0)\|_{\Sigma_0^{-1}}^2 \le 1\right\}\right] \cdot \sqrt{\mathrm{err}_1(\Sigma_0)}$$

$$\le K\sqrt{\varepsilon_{\mathrm{TV}}}$$

The first inequality is Holder's inequality, while the second is an importance weighting argument to replace $a_1 \sim \pi(x_1)$ with the uniform distribution. Next we re-write the expectation using the low rank dynamics, and also use the fact that $x_1 \ne x_{\mathrm{absorb}}$ implies that the previous transition was non-absorbing, which yields

the indicator. Finally, we use the Cauchy-Schwarz inequality in the $\Sigma_0$ norm, along with the implication of the indicator and the assumed bound on $\mathrm{err}_1(\Sigma_0)$. This argument applies as is to all indices $h' > 0$ and for $h' = 0$ we simply apply Holder's inequality and the definition of $\mathrm{err}_0$ to obtain the upper bound $\sqrt{\varepsilon_{\mathrm{TV}}}$. In total, these terms account for the $HK\sqrt{\varepsilon_{\mathrm{TV}}}$ terms on both sides of the lemma statement.

Next we turn to $\mathrm{Term2}_{h'}$. For the first inequality in the lemma statement, we need to upper bound $-\mathrm{Term2}_{h'}$, but this term is easily seen to be non-positive, since $\hat{V}(x_{\mathrm{absorb}}) = 0$ always. So this proves the first inequality. For the second inequality, we have (again focusing on time 1)

$$\mathrm{Term2}_1 = \mathbb{E}\left[\int T_1(x_2 \mid x_1, a_1)\mathbf{1}\{(x_1, a_1) \notin \mathcal{K}_1\}\hat{V}_2(x_2)\right] \leq \mathbb{P}\left[(x_1, a_1) \notin \mathcal{K}_1 \mid \pi, M_\mathcal{K}\right].$$

The same argument applies for all $h \geq 1$. $\qquad\square$

In the next lemma, we consider the case where $\rho_j$ has large escaping probability, measured with respect to the known sets $\mathcal{K}_h(\Sigma_{h,j})$. Recall that $\Sigma_{h,j}$ is the second moment matrix of the true features $\phi_h^\star$ at time $h$ induced by the previous roll-in policies $\rho_0, \ldots, \rho_{j-1}$.

**Lemma 7.** *Consider iteration $j$ of* FLAMBE *and assume that* $\mathrm{err}_h(\Sigma_{h-1,j}) \leq \varepsilon_{\mathrm{TV}}$ *for each $h$ with our current model* $\widehat{\mathcal{M}}$. *Define* $R_h(x, a) := \mathbf{1}\{(x, a) \notin \mathcal{K}_h(\Sigma_{h,j})\}$ *for each $h$. Then,*

$$\max_h \mathrm{tr}\left(\mathbb{E}_{\rho_j}\left[\phi_h^\star(x_h, a_h)\phi_h^\star(x_h, a_h)\right]\Sigma_{h,j-1}^{-1}\right) \geq \frac{1}{H}\max_h\left\{\hat{\mathbb{E}}_{\rho_j}\left[R_h(x_h, a_h)\right] - HK\sqrt{\varepsilon_{\mathrm{TV}}}\right\}$$

*Proof.* For shorthand let $\mathcal{K}_h := \mathcal{K}_h(\Sigma_{h,j})$ denote the known set at round $j$, and let $\mathcal{M}_\mathcal{K}$ denote the corresponding absorbing MDP. Then, applying the second inequality in [Lemma 6](#) we have

$$\hat{\mathbb{E}}_{\rho_j}\left[R_h(x_h, a_h)\right] \leq \mathbb{E}\left[R_h(x_h, a_h) \mid \rho_j, \mathcal{M}_\mathcal{K}\right] + HK\sqrt{\varepsilon_{\mathrm{TV}}} + \sum_{h'=0}^{h-1}\mathbb{P}\left[(x_{h'}, a_{h'}) \notin \mathcal{K}_{h'} \mid \rho_j, \mathcal{M}_\mathcal{K}\right]$$

$$\leq HK\sqrt{\varepsilon_{\mathrm{TV}}} + \sum_{h'=0}^{h}\mathbb{P}\left[(x_{h'}, a_{h'}) \notin \mathcal{K}_{h'} \mid \rho_j, \mathcal{M}_\mathcal{K}\right]$$

$$\leq HK\sqrt{\varepsilon_{\mathrm{TV}}} + \sum_{h'=0}^{h}\mathbb{P}\left[(x_{h'}, a_{h'}) \notin \mathcal{K}_{h'} \mid \rho_j, \mathcal{M}\right]$$

$$= HK\sqrt{\varepsilon_{\mathrm{TV}}} + \sum_{h'=0}^{h}\mathbb{P}\left[\|\phi_{h'}^\star(x_{h'}, a_{h'})\|_{\Sigma_{h',j}^{-1}} \geq 1 \mid \rho_j, \mathcal{M}\right]$$

$$\leq HK\sqrt{\varepsilon_{\mathrm{TV}}} + \sum_{h'=0}^{h}\mathbb{E}_{\rho_j}\left[\mathrm{tr}\left(\phi_{h'}^\star(x_{h'}, a_{h'})\phi_{h'}^\star(x_{h'}, a_{h'})^\top\Sigma_{h',j}^{-1}\right)\right]$$

The last step follows from Markov's inequality. Since both matrices are positive semidefinite, the trace terms are all non-negative. Therefore, by the pigeonhole principle, there exists some $h' \in \{0, \ldots, h\}$ for which

$$\mathbb{E}_{\rho_j}\left[\mathrm{tr}\left(\phi_{h'}^\star(x_{h'}, a_{h'})\phi_{h'}^\star(x_{h'}, a_{h'})^\top\Sigma_{h',j}^{-1}\right)\right] \geq \frac{1}{H}\left(\hat{\mathbb{E}}_{\rho_j}\left[R_h(x_h, a_h)\right] - HK\sqrt{\varepsilon_{\mathrm{TV}}}\right).$$

This argument applies for all $R_h$, and so we obtain the lemma. $\qquad\square$

Next we argue that there cannot be too many iterations for which $\max_h \hat{\mathbb{E}}_{\rho_j}\left[R_h(x_h, a_h)\right]$ is large. For notation, here we use $\widehat{\mathcal{M}}^{(j)}$ to denote the MDP model in iteration $j$ and we use $R_h^{(j)}$ to denote the reward functions in [Lemma 7](#) derived from the known sets in iteration $j$.

20

**Corollary 8.** *Assume that for each round $j \in [J_{\max}]$ and for all $h$ we have $\mathrm{err}_h(\Sigma_{h-1,j}) \leq \varepsilon_{\mathrm{TV}}$. Set*

$$J_{\max} := \frac{4Hd}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}} \cdot \log\left(1 + \frac{4H}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}}\right). \tag{4}$$

*Then there exists some $j \in [J_{\max}]$ for which $\max_h \mathbb{E}\left[R_h^{(j)}(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] \leq 2HK\sqrt{\varepsilon_{\mathrm{TV}}}$.*

*Proof.* Suppose that in round $j$, it holds that $\max_h \mathbb{E}\left[R_h^{(j)}(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] \geq 2HK\sqrt{\varepsilon_{\mathrm{TV}}}$. Then, by Lemma 7, there exists some time step $h$ for which

$$\mathrm{tr}\left(\mathbb{E}_{\rho_j}\left[\phi_h^\star(x_h, a_h)\phi_h^\star(x_h, a_h)^\top\right] \Sigma_{h,j}^{-1}\right) \geq K\sqrt{\varepsilon_{\mathrm{TV}}}.$$

Note that we also have $\Sigma_{h,j+1} = \Sigma_{h,j} + \mathbb{E}_{\rho_j}\left[\phi_h^\star(x_h, a_h)\phi_h^\star(x_h, a_h)^\top\right]$, so we are in a position to apply the elliptical potential argument. Specifically if $J$ is the number of iterations for which the above inequality holds for some $h$, then applying Lemma 26 for each $h$ and summing across $h$ yields

$$JK\sqrt{\varepsilon_{\mathrm{TV}}} \leq (1 + 1/\lambda)dH \log(1 + J_{\max}/d)$$

Plugging in our choice of $J_{\max}$, and using the fact that $\lambda < 1$ we have

$$J < \frac{2Hd}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}} \log\left(1 + \frac{4H}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}} \log\left(1 + \frac{4H}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}}\right)\right)$$

$$\leq \frac{2Hd}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}} \log\left(1 + \left(\frac{4H}{\lambda K \sqrt{\varepsilon_{\mathrm{TV}}}}\right)^2\right) \leq J_{\max}.$$

This means that in $J_{\max}$ iterations, we can have $\max_h \mathbb{E}\left[R_h^{(j)}(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] \geq 2HK\sqrt{\varepsilon_{\mathrm{TV}}}$ in at most $J < J_{\max}$ of them. Thus we must have one where this quantity is small, which proves the lemma. $\qquad\square$

Next, we state a guarantee provided by Algorithm 2, which is a more convenient form of Lemma 17.

**Lemma 9.** *Fix any iteration $j$, time $h$, function $f : \mathcal{X} \times \mathcal{A} \to [0, 1]$, policy $\pi$, any $\alpha > 0$. Then*

$$\mathbb{E}\left[f(x_h, a_h) \mid \pi, \widehat{\mathcal{M}}^{(j)}\right] \leq \frac{T\beta}{2\alpha} + \frac{\alpha d}{2T} + \frac{\alpha KH}{2}\mathbb{E}\left[f(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right],$$

*where $T \leq 4d\log(1 + 4/\beta)/\beta$ and $\beta > 0$ is the parameter to Algorithm 2.*

*Proof.* We suppress the dependence on $j$. Let us first focus on $\rho_h^{\mathrm{pre}}$, which is output of Algorithm 2 for some time step $h$. $\rho_h^{\mathrm{pre}}$ induces a distribution over states at time step $h$, and we argue that this distribution adequately covers all possible roll-in distributions in the model $\widehat{\mathcal{M}} = \widehat{\mathcal{M}}^{(j)}$. Consider any function $f : \mathcal{X} \times \mathcal{A} \to [0, 1]$, any policy $\pi$, any $\Sigma \succ 0$, and $\alpha > 0$. Calling $f_\pi(x_h) = \int \pi(a_h \mid x_h)f(x_h)$, we have

$$\hat{\mathbb{E}}_\pi f(x_h, a_h) = \hat{\mathbb{E}}_\pi \left\langle \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}), \int \hat{\mu}_{h-1}(x_h)f_\pi(x_h) \right\rangle$$

$$\leq \hat{\mathbb{E}}_\pi \left\| \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}) \right\|_{\Sigma^{-1}} \cdot \left\| \int \hat{\mu}_{h-1}(x_h)f_\pi(x_h) \right\|_{\Sigma}$$

$$\leq \frac{1}{2\alpha}\hat{\mathbb{E}}_\pi \|\phi_{h-1}(x_{h-1}, a_{h-1})\|_{\Sigma^{-1}}^2 + \frac{\alpha}{2}\left\| \int \hat{\mu}_{h-1}(x_h)f_\pi(x_h) \right\|_{\Sigma}^2.$$

21

Here we expand $\hat{T}_{h-1}$ in terms of its low rank representation and then apply the Cauchy-Schwarz inequality in the norm induced by $\Sigma$. Finally we use the AM-GM inequality which holds for any non-negative $\alpha$.

We instantiate $\Sigma$ to be the covariance matrix induced by $\rho_h^{\mathrm{pre}}$. First, for any policy $\pi$ we define the $h-1$ step model covariance as $\Sigma_\pi := \hat{\mathbb{E}}_\pi \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}) \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})^\top$, where the dependence on $h-1$ is suppressed in the notation. Note that both the expectation and the embedding are taken with respect to the model $\widehat{\mathcal{M}}$. Then, the output of Algorithm 2 is a $h$-step policy $\rho_h^{\mathrm{pre}}$ that is defined as a mixture over $T$ policies $\pi_1, \ldots, \pi_T$. Using these policies, we define $\Sigma$ as follows:

$$\Sigma = \Sigma_{\rho_h^{\mathrm{pre}}} + \frac{I_{d\times d}}{T} = \frac{1}{T}\sum_{t=1}^T \Sigma_{\pi_t} + \frac{I_{d\times d}}{T}.$$

As we run Algorithm 2 using $\hat{T}_{0:h-1}$ we can apply Lemma 17 on the $h$ step MDP $\hat{T}_{0:h-1}$. In other words, in Lemma 17, we set $H \leftarrow h$ and $\widetilde{\mathcal{M}} \leftarrow \widehat{\mathcal{M}}$. The conclusion is that $T \le 4d\log(1 + 4/\beta)/\beta$, where $\beta$ is the parameter to the subroutine, and we can also bound the first term above:

$$\hat{\mathbb{E}}_\pi \left\| \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}) \right\|_{\Sigma^{-1}}^2 = \hat{\mathbb{E}}_\pi \phi_{h-1}(x_{h-1}, a_{h-1})^\top \left( \Sigma_{\rho_h^{\mathrm{pre}}} + \frac{I_{d\times d}}{T} \right)^{-1} \phi_{h-1}(x_{h-1}, a_{h-1}) \le T\beta.$$

Next, we turn to the second term. Expanding the definition of $\Sigma$, we have

$$\left\| \int \hat{\mu}_{h-1}(x_h) f_\pi(x_h) \right\|_\Sigma^2$$

$$= \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}}} \left( \left\langle \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}), \int \hat{\mu}_{h-1}(x_h) f_\pi(x_h) \right\rangle \right)^2 + \frac{\left\| \int \hat{\mu}_{h-1}(x_h) f_\pi(x_h) \right\|_2^2}{T}$$

$$= \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}}} \left( \hat{\mathbb{E}}\left[ f_\pi(x_h) \mid x_{h-1}, a_{h-1} \right] \right)^2 + \frac{\left\| \int \hat{\mu}_{h-1}(x_h) f_\pi(x_h) \right\|_2^2}{T}$$

$$\le \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}}} f_\pi(x_h) + \frac{\left\| \int \hat{\mu}_{h-1}(x_h) f_\pi(x_h) \right\|_2^2}{T} \le \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}}} f_\pi(x_h) + d/T.$$

The first inequality is Jensen's inequality along with the fact that $f(x_h)^2 \le f(x_h)$ since $f : \mathcal{X} \to [0, 1]$. The second inequality is based on our normalization assumptions on $\mu_{h-1}$, which we also impose on $\hat{\mu}_{h-1}$. Finally, collecting all the terms and importance weighting the last action, we obtain the bound

$$\hat{\mathbb{E}}_\pi f(x_h) \le \frac{T\beta}{2\alpha} + \frac{\alpha K}{2} \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}} \circ \mathrm{unif}(\mathcal{A})} f(x_h, a_h) + \frac{\alpha d}{2T}.$$

This bound applies to $\rho_h^{\mathrm{pre}}$. As $\rho_j$ is a uniform mixture of these policies and as $f$ is non-negative, we see that $\hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}} \circ \mathrm{unif}(\mathcal{A})} f(x_h, a_h) \le H \cdot \hat{\mathbb{E}}_{\rho_j} f(x_h, a_h)$, which proves the lemma. $\qquad\square$

Finally, we use the guarantee for Algorithm 2, to prove that our model $\widehat{\mathcal{M}}$ universally approximate the true MDP as soon as $\max_h \mathbb{E}\left[ R_h^{(j)}(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)} \right] \le 2HK\sqrt{\varepsilon_{\mathrm{TV}}}$. For the lemma, we use the concept of a sparse reward function. $R : \mathcal{X} \times \mathcal{A} \to [0, 1]$ is called *sparse* if all value functions are in $[0, 1]$. For example, this holds if $R$ is only associated with state-action pairs at a single time point.

**Lemma 10.** *Assume that for each round $j \in [J_{\max}]$ and for all $h$, we have $\mathrm{err}_h(\Sigma_{h-1,j}) \le \varepsilon_{\mathrm{TV}}$, and set $J_{\max}$ as in (4). Then the final MDP model $\widehat{\mathcal{M}}$ satisfies the following guarantee: For any sparse reward function $R : \mathcal{X} \times \mathcal{A} \to [0, 1]$, any policy $\pi$, and any $\alpha > 0$ we have*

$$\left| V(\pi; R, \widehat{\mathcal{M}}) - V(\pi; R, \mathcal{M}) \right| \le HK\sqrt{\varepsilon_{\mathrm{TV}}} + H\varepsilon_{\mathrm{escape}},$$

*where* $\varepsilon_{\text{escape}} := \alpha H^2 K^2 \sqrt{\varepsilon_{\text{TV}}} + \frac{T\beta}{2\alpha} + \frac{\alpha d}{2T} + HK\sqrt{\varepsilon_{\text{TV}}}$, $T \le 4d\log(1+4/\beta)/\beta$ *and* $\beta > 0$ *is the parameter to* Algorithm 2.

*Proof.* Via Corollary 8, there must be some round $j$ for which

$$\max_h \mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] = \max_h \mathbb{E}\left[R_h^{(j)}(x_h, a_h) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] \le 2HK\sqrt{\varepsilon_{\text{TV}}}.$$

We will prove the guarantee for this round $j$, and at the end of the proof argue that this also applies to the final learned model.

Combining the lower bound of the simulation lemma (Lemma 6) with the planning guarantee (Lemma 9) we see that, for any policy $\pi$

$$
\begin{aligned}
\mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \pi, \mathcal{M}_{\mathcal{K}}^{(j)}\right] &\le \mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j} \mid \pi, \widehat{\mathcal{M}}^{(j)}\right] + HK\sqrt{\varepsilon_{\text{TV}}} \\
&\le \frac{\alpha K H}{2}\mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \rho_j, \widehat{\mathcal{M}}^{(j)}\right] + \frac{T\beta}{2\alpha} + \frac{\alpha d}{2T} + HK\sqrt{\varepsilon_{\text{TV}}} \\
&\le \alpha K^2 H^2 \sqrt{\varepsilon_{\text{TV}}} + \frac{T\beta}{2\alpha} + \frac{\alpha d}{2T} + HK\sqrt{\varepsilon_{\text{TV}}} =: \varepsilon_{\text{escape}}
\end{aligned}
$$

Now that we have upper bounded the escaping probability, we can turn to the approximation guarantee. While we are not in the exact setting of Lemma 6, since we have a sparse reward function, all values are in $[0, 1]$ so the same argument applies. For one side of the error guarantee, since we assume that $\text{err}_h(\Sigma_{h-1,j}) \le \varepsilon_{\text{TV}}$ for all iterations, we have

$$
\begin{aligned}
V(\pi; R, \widehat{\mathcal{M}}^{(j)}) &\le V(\pi; R, \mathcal{M}_{\mathcal{K}}^{(j)}) + HK\sqrt{\varepsilon_{\text{TV}}} + \sum_{h=0}^{H-1} \mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \pi, \mathcal{M}_{\mathcal{K}}^{(j)}\right] \\
&\le V(\pi; R, \mathcal{M}) + HK\sqrt{\varepsilon_{\text{TV}}} + \sum_{h=0}^{H-1} \mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \pi, \mathcal{M}_{\mathcal{K}}^{(j)}\right] \\
&\le V(\pi; R, \mathcal{M}) + HK\sqrt{\varepsilon_{\text{TV}}} + H\varepsilon_{\text{escape}}.
\end{aligned}
$$

Here the first inequality is Lemma 6, while the second is due to (3). For the other direction, we first use (3) and then Lemma 6:

$$
\begin{aligned}
V(\pi; R, \mathcal{M}) &\le V(\pi; R, \mathcal{M}_{\mathcal{K}}^{(j)}) + \sum_{h=0}^{H-1} \mathbb{P}\left[(x_h, a_h) \notin \mathcal{K}_h(\Sigma_{h,j}) \mid \pi, \mathcal{M}_{\mathcal{K}}^{(j)}\right] \\
&\le V(\pi; R, \widehat{\mathcal{M}}^{(j)}) + HK\sqrt{\varepsilon_{\text{TV}}} + H\varepsilon_{\text{escape}}.
\end{aligned}
$$

This proves the result for the MDP model $\widehat{\mathcal{M}}^{(j^\star)}$ at the time $j^\star$ where the exploratory policy $\rho_{j^\star}$ fails to achieve large reward on $R_h^{(j^\star)}$. We now claim this applies for all iterations after $j^\star$, and in particular it holds at the end of the algorithm. To see why, observe that $\Sigma_{h,j+1} \succeq \Sigma_{h,j}$ for all $j$, and so $\mathcal{K}_h(\Sigma_{h,j}) \subset \mathcal{K}_h(\Sigma_{h,j+1})$ for all rounds. Since the known set increases with $j$, the escaping probability is decreasing, so it is upper bounded by $\varepsilon_{\text{escape}}$ for all rounds after $j \ge j^\star$. Additionally, we assume that $\text{err}_h(\Sigma_{h-1,j}) \le \varepsilon_{\text{TV}}$ for all $j \in [J_{\max}]$, so the total variation term in Lemma 6 remains bounded. Working through the proof of Lemma 6, we can see that the bound continues to hold for all $j^\star \le j \le J_{\max}$, which proves the result. $\qquad\square$

**Final steps.** Let us collect all of the conditions and bounds here. At the end of the algorithm, we have

$$\max_{\pi, R}\left|V(\pi; R, \widehat{\mathcal{M}}) - V(\pi; R, \mathcal{M})\right| \le HK\sqrt{\varepsilon_{\text{TV}}} + H\varepsilon_{\text{escape}}, \tag{5}$$

where we may set

$$\varepsilon_{\text{escape}} := \min_{\alpha > 0} \left\{ \alpha H^2 K^2 \sqrt{\varepsilon_{\text{TV}}} + \frac{T\beta}{2\alpha} + \frac{\alpha d}{2T} + HK\sqrt{\varepsilon_{\text{TV}}} \right\}, \qquad T \leq \frac{4d \log(1 + 4/\beta)}{\beta},$$

and $\beta > 0$ is the parameter to [Algorithm 2]. Applying [Corollary 5] and taking a union bound over all iterations $j \in [J_{\max}]$ and all times $h$, we can set

$$\varepsilon_{\text{TV}} := \lambda d + \frac{2 \log(J_{\max} H |\Phi| |\Upsilon| / \delta)}{n},$$

where $\lambda > 0$ is a parameter in the analysis. Finally, the total number of samples collected is

$$nH J_{\max}, \qquad \text{where,} \qquad J_{\max} := \frac{4Hd}{\lambda K \sqrt{\varepsilon_{\text{TV}}}} \cdot \log\left(1 + \frac{4H}{\lambda K \sqrt{\varepsilon_{\text{TV}}}}\right)$$

We start by optimizing for $\alpha$ in the definition of $\varepsilon_{\text{escape}}$, which yields $\alpha = \sqrt{\frac{T\beta}{2(H^2 K^2 \sqrt{\varepsilon_{\text{TV}}} + d/(2T))}}$. Plugging into $\varepsilon_{\text{escape}}$ and using the bound on $T$, we get

$$\varepsilon_{\text{escape}} \leq \sqrt{2T\beta} \cdot \left(HK\varepsilon_{\text{TV}}^{1/4} + \sqrt{d/(2T)}\right) + HK\sqrt{\varepsilon_{\text{TV}}}$$
$$\leq 2\sqrt{8d \log(1 + 4/\beta)} HK\varepsilon_{\text{TV}}^{1/4} + \sqrt{d\beta}.$$

Here we are using the fact that $\varepsilon_{\text{TV}} \leq 1$, which is without loss of generality, since [(5)] is trivial when $\varepsilon_{\text{TV}} \geq 1$. Now we set $\beta = H^2 K^2 \sqrt{\varepsilon_{\text{TV}}}$ so that

$$\varepsilon_{\text{escape}} \leq 16\sqrt{d \log(1 + 4/\varepsilon_{\text{TV}})} HK\varepsilon_{\text{TV}}^{1/4}.$$

Thus, we may restate the final accuracy guarantee as

$$\max_{\pi, R} \left| V(\pi; R, \widehat{\mathcal{M}}) - V(\pi; R, \mathcal{M}) \right| \leq HK\sqrt{\varepsilon_{\text{TV}}} + 16\sqrt{d \log(1 + 4/\varepsilon_{\text{TV}})} H^2 K\varepsilon_{\text{TV}}^{1/4}$$
$$\leq 17\sqrt{d \log(1 + 4/\varepsilon_{\text{TV}})} H^2 K\varepsilon_{\text{TV}}^{1/4}.$$

We want this to be upper bounded by $\varepsilon$, the final accuracy parameter, which means we can take

$$\varepsilon_{\text{TV}} = c\frac{\varepsilon^4 H^{-8} K^{-4} d^{-2}}{\log^2(1 + 1/\varepsilon)},$$

where $c > 0$ is a universal constant. Looking at the definition of $\varepsilon_{\text{TV}}$ and $T_{\max}$, we set

$$\lambda = c\frac{\varepsilon^4 H^{-8} K^{-4} d^{-3}}{\log^2(1 + 1/\varepsilon)}, \quad T_{\max} = \tilde{O}\left(\frac{H^{13} d^5 K^5}{\varepsilon^6}\right), \quad n = \tilde{O}\left(\frac{H^8 K^4 d^2 \log(|\Phi| |\Upsilon| / \delta)}{\varepsilon^4}\right),$$

where we are ignoring logarithmic factors. This gives the final sample complexity of

$$\tilde{O}\left(\frac{H^{22} K^9 d^7 \log(|\Phi| |\Upsilon| / \delta)}{\varepsilon^{10}}\right).$$

Finally, note that [(1)] is implied by the final accuracy guarantee, since we may choose $R$ to be the total variation distance between our model and the true transition dynamics at time $h$, which is clearly a sparse reward function.

**Analysis with the sampling oracle.** With the sampling oracle, the argument is very similar. The main difference is in Lemma 9, which is the only place where we use the exact planner. Instead, we modify the proof of Lemma 9 to instead use Lemma 18 to obtain $\hat{\Sigma}$ and $\rho_h^{\mathrm{pre}}$, and we do the Cauchy-Schwarz step using $\hat{\Sigma}$. By Lemma 18 the first term is still $O(T\beta)$ and for the second term we pay an additive $O(\beta)$ to translate from $\hat{\Sigma}$ to $\Sigma$ (since the spectral norm error is $O(\beta/d)$ and Euclidean norm of the term involving $\hat{\mu}_{h-1}$ is at most $d$). This we have an additional $O(\alpha\beta)$ term in the sample-based analog of Lemma 9. However, as we use the bound $T \leq O(d\log(1 + 1/\beta)/\beta)$ in the remaining calculations, the new $O(\alpha\beta)$ term is only larger than the $O(\alpha d/T)$ term by a logarithmic factor. In particular, above we have a $\sqrt{d\beta}$ term in the bound for $\varepsilon_{\mathrm{escape}}$, but with the sampling oracle, we will additionally have a $O(\sqrt{T}\beta) = O(\sqrt{d\beta\log(1 + 1/\beta)})$ term in this bound. Ultimately, this only affects the final sample complexity bound in logarithmic factors. We adjust the failure probability accordingly, using $\delta/2$ probability for invocations of Lemma 18, and $\delta/2$ for the invocations of Corollary 5. As we invoke the planner polynomially many times, the total number of calls to the sampling oracle is polynomial in all parameters.

## B.1 Refined analysis for simplex representations.

Here we prove Theorem 3 by considering a different potential function argument and a different instantiation of the planning algorithm that directly attempts to visit each latent state. In particular, we instantiate Algorithm 1 with the planning routine presented in Algorithm 3. Note that this planner does not require the parameter $\beta$, but it does assume that $\hat{\phi}(x, a) \in \Delta([d_{\mathrm{LV}}])$ for each $(x, a)$.

The analog of $\Sigma_{h,j}$ is the cumulative probability of hitting latent variables $z_h \in \mathcal{Z}_h$. Formally, we define

$$p_{h,j}(z) := \sum_{i=0}^{j-1} \mathbb{P}\left[z_h = z \mid \rho_i, \mathcal{M}\right].$$

We make two remarks. First $p_{h,j}$ is not a distribution, rather it is the sum of $j$ probability distributions. Second, we use $p_{h,j}$ to measure the coverage at time $h$, since $z_h$ is the latent variable that generates $x_h$. This indexing is different from how we use $\Sigma_{h,j}$ to measure coverage at time $h + 1$ in the general case.

We now state the analog of Corollary 5.

**Corollary 11.** *For $j \geq 1, h \in \{0, \ldots, H-1\}, \delta \in (0, 1)$ and let $\rho_0, \ldots, \rho_{j-1}$ be any (possibly data-dependent) policies, with $p_{h,j}$ defined accordingly. Let $D_h$ be a dataset of $nj$ examples, where for each $0 \leq i < j$, we collect $n$ triples $(x_h, a_h, x_{h+1})$ by rolling in with $\rho_i$ to $x_h$ and taking $a_h$ uniformly at random. Then, with probability at least $1 - \delta$ the output $(\hat{\phi}_h, \hat{\mu}_h)$ of $\mathrm{MLE}(D_h)$ satisfies*

$$\sum_{z \in \mathcal{Z}_h} p_{h,j}(z) \mathbb{E}\left[\left\|\left\langle \hat{\phi}(x_h, a_h), \hat{\mu}_h(\cdot)\right\rangle - T_h(\cdot \mid x_h, a_h)\right\|_{\mathrm{TV}}^2 \mid z_h = z, a_h \sim \mathrm{unif}(\mathcal{A})\right] \leq \frac{2\log(|\Phi||\Upsilon|/\delta)}{n}.$$

*Proof.* This is an immediate consequence of Theorem 21, using the definition of $p_{h,j}$. $\square$

We denote the LHS of the above lemma as $\mathrm{err}_h(p_{h,j})$. Now we define the known set $\mathcal{K}_h$ and the absorbing MDP. In the simplex features setting, the known set $\mathcal{K}_h$ is instead defined in terms of latent variables. Recall that we can augment every trajectory $\tau$ with the latent variables generated along the trajectory that is $\tau = (z_0, x_0, a_0, z_1, x_1, a_1, \ldots, z_{H-1}, x_{H-1}, a_{H-1})$. We therefore define $\mathcal{K}_{h,j} := \{z \in \mathcal{Z}_h : p_{h,j}(z) \geq \Delta\}$ where $\Delta$ is some parameter we will set towards the end of the proof. The absorbing MDP $\mathcal{M}_{\mathcal{K}}$ in iteration $j$ is defined to have transition operator that, for each $h$, transitions from $z_h$ to $x_{\mathrm{absorb}}$ if $z_h \notin \mathcal{K}_{h,j}$ and otherwise transitions as in $\mathcal{M}$. As in the more general analysis, $x_{\mathrm{absorb}}$ is an absorbing state with a single self-looping action $a_{\mathrm{absorb}}$ and we always consider $(x_{\mathrm{absorb}}, a_{\mathrm{absorb}})$ to be known.

We now state the analog of Lemma 6.

**Lemma 12.** *Let $\hat{\phi}_{0:H-1}, \hat{\mu}_{0:H-1}$ be an MDP model with simplex features and let $p_{0:H-1}$ be non-negative vectors. Assume that $\mathrm{err}_h(p_{h-1}) \leq \varepsilon_{\mathrm{TV}}$ for each $h$. Let $f : \mathcal{X} \times \mathcal{A} \to [0,1]$ be any function such that $f(x_{\mathrm{absorb}}, a_{\mathrm{absorb}}) = 0$ and let $\pi$ be any policy. Then, for any $h$*

$$\mathbb{E}_\pi\left[f(x_h, a_h) \mid \mathcal{M}_\mathcal{K}\right] - HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} \leq \hat{\mathbb{E}}_\pi\left[f(x_h, a_h)\right] \leq \mathbb{E}_\pi\left[f(x_h, a_h) \mid \mathcal{M}_\mathcal{K}\right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$
$$+ \sum_{h'=1}^{h} \mathbb{P}\left[z_{h'} \notin \mathcal{K}_{h'} \mid \pi, \mathcal{M}_\mathcal{K}\right]$$

*Proof.* As in the proof of [Lemma 6](#), we must control two terms for each $h' < h$:

$$\mathrm{Term1}_{h'} := \mathbb{E}\left[\int (\hat{T}_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}) - T_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}))\hat{V}_{h'+1}(x_{h'+1}) \mid \pi, \mathcal{M}_\mathcal{K}\right]$$

$$\mathrm{Term2}_{h'} := \mathbb{E}\left[\int (T_{h'}(x_{h'+1} \mid x_{h'}, a_{h'}) - T_{h',\mathcal{K}}(x_{h'+1} \mid x_{h'}, a_{h'}))\hat{V}_{h'+1}(x_{h'+1}) \mid \pi, \mathcal{M}_\mathcal{K}\right].$$

For $\mathrm{Term1}_{h'}$, as we are in $\mathcal{M}_\mathcal{K}$ we can ignore the trajectories where $z_{h'} \notin \mathcal{K}_{h'}$. Thus considering $h' = 1$

$$\mathrm{Term1}_1 = \sum_{z \in \mathcal{K}_1} \mathbb{P}_\pi[z_1 = z \mid \mathcal{M}_\mathcal{K}] \cdot \mathbb{E}_\pi\left[\int (\hat{T}_1(x_2 \mid x_1, a_1) - T_1(x_2 \mid x_1, a_1))\hat{V}_2(x_2) \mid z_1 = z\right]$$

$$\leq K \sum_{z \in \mathcal{K}_1} \mathbb{P}_\pi[z_1 = z \mid \mathcal{M}_\mathcal{K}] \cdot \mathbb{E}\left[\left\|\hat{T}_1(\cdot \mid x_1, a_1) - T_1(\cdot \mid x_1, a_1)\right\|_{\mathrm{TV}} \mid z_1 = z, a_1 \sim \mathrm{unif}(\mathcal{A})\right]$$

$$\leq K \sqrt{\sum_{z \in \mathcal{K}_1} \mathbb{P}_\pi[z_1 = z \mid \mathcal{M}_\mathcal{K}] \cdot \mathbb{E}\left[\left\|\hat{T}_1(\cdot \mid x_1, a_1) - T_1(\cdot \mid x_1, a_1)\right\|_{\mathrm{TV}}^2 \mid z_1 = z, a_1 \sim \mathrm{unif}(\mathcal{A})\right]}$$

$$\leq K\sqrt{\mathrm{err}_1(p_1)/\Delta} \leq K\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}.$$

Here we are using that the total variation term is non-negative, and that $\mathbb{P}_\pi[z_1 = z \mid \mathcal{M}_\mathcal{K}] \leq 1$, while $p_1(z) \geq \Delta$ by the fact that $z \in \mathcal{K}_1$. This argument applies for all $h'$ and yields the $HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$ term on both sides of the statement. For $\mathrm{Term2}_{h'}$, we clearly have

$$\mathrm{Term2}_{h'} \leq \mathbb{P}\left[z_{h'+1} \notin \mathcal{K}_{h'+1} \mid \pi, \mathcal{M}_\mathcal{K}\right].$$

As in the proof of [Lemma 6](#), $\mathrm{Term2}_{h'} \geq 0$ which yields the lower bound. $\qquad\square$

Next we argue that if the exploratory policy $\rho_j$ that we find has large escaping probability then we will add some latent variable to the known set in the next iteration.

**Lemma 13.** *Consider iteration $j$ of FLAMBE and assume that $\mathrm{err}_h(p_{h,j}) \leq \varepsilon_{\mathrm{TV}}$ for each $h$. Define $R_h(x, a) := \sum_{z \notin \mathcal{K}_{h,j}} \phi_h^\star(x, a)[z]$. Then*

$$\max_h \mathbb{P}\left[z_h \notin \mathcal{K}_{h,j} \mid \rho_j\right] \geq \frac{1}{H} \max_h \left\{\hat{\mathbb{E}}_{\rho_j}[R_h(x_h, a_h)] - HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}\right\}.$$

*In particular, if there exists some $h$ such that $\hat{\mathbb{E}}_{\rho_j}[R_h(x_h, a_h)] \geq HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + Hd_{\mathrm{LV}}\Delta$, then there exists some $h', z \notin \mathcal{K}_{h'}$ such that $\mathbb{P}\left[z_{h'} = z \mid \rho_j\right] \geq \Delta$.*

*Proof.* Observe that by the definition of $R_h$, we have

$$\hat{\mathbb{E}}_{\rho_j}\left[R_h(x_h, a_h)\right] \leq \mathbb{E}_{\rho_j}[R_h(x_h, a_h) \mid \mathcal{M}_\mathcal{K}] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + \sum_{h'=1}^{h} \mathbb{P}\left[z_{h'} \notin \mathcal{K}_{h',j} \mid \rho_j, \mathcal{M}_\mathcal{K}\right]$$

$$= HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + \sum_{h'=1}^{h+1} \mathbb{P}\left[z_{h'} \notin \mathcal{K}_{h',j} \mid \rho_j, \mathcal{M}_\mathcal{K}\right].$$

Both statements now follow from the pigeonhole principle. $\qquad\square$

Next we prove the analog of Lemma 10. For this, we compute $\rho_h^{\mathrm{pre}}$ using the planning routine in Algorithm 3, with the planning guarantee in Lemma 16.

**Lemma 14.** *Assume that for each round $j \in [J_{\max}]$ and for all $h$, we have $\mathrm{err}_h(p_{h,j}) \le \varepsilon_{\mathrm{TV}}$ and set $J_{\max} = Hd_{\mathrm{LV}} + 1$. Then the final MDP model $\widehat{\mathcal{M}}$ satisfies the following guarantee: For any sparse reward function $R : \mathcal{X} \times \mathcal{A} \to [0,1]$ and any policy $\pi$, we have*

$$\left| V(\pi; R, \widehat{\mathcal{M}}) - V(\pi; R, \mathcal{M}) \right| \le HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + H\varepsilon_{\mathrm{escape}},$$

*where $\varepsilon_{\mathrm{escape}} := H^2 K d_{\mathrm{LV}}^2 \Delta + \left( H^2 K^2 d_{\mathrm{LV}} + HK d_{\mathrm{LV}} \right) \sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$.*

*Proof.* First observe that by Lemma 13, in every iteration where $\rho_j$ satisfies $\max_h \hat{\mathbb{E}}_{\rho_j} [R_h(x_h, a_h)] \ge HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + Hd_{\mathrm{LV}}\Delta$, we add some latent variable at some time point to the known set. This means that this can only happen for at most $Hd_{\mathrm{LV}}$ iterations, and so, by the setting of $J_{\max}$, there must be some iteration $j$ in which $\max_h \hat{\mathbb{E}}_{\rho_j} [R_h(x_h, a_h)] \le HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + Hd_{\mathrm{LV}}\Delta$. In this iteration $j$, we have

$$\mathbb{P}\left[ z_{h+1} \notin \mathcal{K}_{h+1,j} \mid \pi, \mathcal{M}_{\mathcal{K}}^{(j)} \right] \le \mathbb{P}\left[ z_{h+1} \notin \mathcal{K}_{h+1,j} \mid \pi, \widehat{\mathcal{M}}^{(j)} \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$

$$= \sum_{z \in \hat{\mathcal{Z}}_h} \hat{\mathbb{E}}_\pi \left[ \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[z] \right] \cdot \mathbb{P}\left[ \bar{\mathcal{K}}_{h+1,j} \mid \pi, \widehat{\mathcal{M}}^{(j)}, \hat{z}_h = i \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$

$$\le K \sum_{z \in \hat{\mathcal{Z}}_h} \hat{\mathbb{E}}_\pi \left[ \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[z] \right] \cdot \mathbb{P}\left[ \bar{\mathcal{K}}_{h+1,j} \mid \widehat{\mathcal{M}}^{(j)}, \hat{z}_h = z, a_h \sim \mathrm{unif}(\mathcal{A}) \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$

$$\le K d_{\mathrm{LV}} \sum_{z \in \hat{\mathcal{Z}}_h} \hat{\mathbb{E}}_{\rho_h^{\mathrm{pre}}} \left[ \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[z] \right] \cdot \mathbb{P}\left[ \bar{\mathcal{K}}_{h+1,j} \mid \widehat{\mathcal{M}}^{(j)}, \hat{z}_h = z, a_h \sim \mathrm{unif}(\mathcal{A}) \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$

$$\le HK d_{\mathrm{LV}} \hat{\mathbb{P}}_{\rho_j} \left[ \bar{\mathcal{K}}_{h+1,j} \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} \le HK d_{\mathrm{LV}} \hat{\mathbb{E}}_{\rho_j} \left[ R_h(x_h, a_h) \right] + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$$

$$\le HK d_{\mathrm{LV}} \left( HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + Hd_{\mathrm{LV}}\Delta \right) + HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} =: \varepsilon_{\mathrm{escape}}$$

Here the first inequality is Lemma 12, while the first equality re-writes the expectation in terms of the latent variable $z_h$. In the second inequality we translate to taking $a_h$ uniformly via importance weighting, while in the third, we apply Lemma 16, which lets us translate to $\rho_h^{\mathrm{pre}}$. Finally, we use that $\rho_j$ uses $\rho_h^{\mathrm{pre}}$ with probability $1/H$ and the definition of $R_h$. The result now follows from Lemma 12, along with the analog of (3). As in the general case, this bound applies for all iterations after the first one where the escaping probability for $\rho_j$ is small. $\qquad\square$

**Final steps.** The final steps with simplex features are much more straightforward than in the general case. First we choose $\Delta$ to balance the two terms in $\varepsilon_{\mathrm{escape}}$. We set $\Delta = (K/d_{\mathrm{LV}})^{2/3}\varepsilon_{\mathrm{TV}}^{1/3}$ which yields

$$\varepsilon_{\mathrm{escape}} \le 3H^2 K d_{\mathrm{LV}} \left( K^{2/3}(d_{\mathrm{LV}}\varepsilon_{\mathrm{TV}})^{1/3} \right),$$

where we are also using the fact that $\varepsilon_{\mathrm{TV}} \le 1$. Via Lemma 14, after $J_{\max} = Hd_{\mathrm{LV}} + 1$ iterations, we are guaranteed that

$$\max_{\pi, R} \left| V(\pi; R, \widehat{\mathcal{M}}) - V(\pi; R, \mathcal{M}) \right| \le HK\sqrt{\varepsilon_{\mathrm{TV}}/\Delta} + H\varepsilon_{\mathrm{escape}} \le O\left( H^3 K^{5/3} d_{\mathrm{LV}}^{4/3} \varepsilon_{\mathrm{TV}}^{1/3} \right).$$

For this to be at most $\varepsilon$ we should set $\varepsilon_{\mathrm{TV}} \leq O\left(\varepsilon^3 H^{-9} K^{-5} d_{\mathrm{LV}}^{-4}\right)$. Applying Corollary 11 and taking a union over all $J_{\max}$ rounds, we want to set

$$n = \frac{2\log(J_{\max}|\Phi||\Upsilon|/\delta)}{\varepsilon_{\mathrm{TV}}} = \tilde{O}\left(\frac{H^9 K^5 d_{\mathrm{LV}}^4 \log(|\Phi||\Upsilon|/\delta)}{\varepsilon^3}\right).$$

The total sample complexity is $nHJ_{\max} = \tilde{O}\left(\frac{H^{11} K^5 d_{\mathrm{LV}}^5 \log(|\Phi||\Upsilon|/\delta)}{\varepsilon^3}\right)$. As in the general setting, the value function guarantee implies (1), which yields the result.

**Analysis with a sampling oracle.** With a sampling oracle, the only difference is in the proof of Lemma 14. Here, we can only apply the second statement of Lemma 16, which yields an additive $O(Kd_{\mathrm{LV}}\varepsilon_{\mathrm{opt}})$ term. By taking $\varepsilon_{\mathrm{opt}} = (H/d_{\mathrm{LV}})\sqrt{\varepsilon_{\mathrm{TV}}/\Delta}$, this additional term can be absorbed into the other additive term at the expense of a constant. Thus we obtain the same guarantee, up to constants, with polynomially many calls to the sampling oracle.

## C   Planning Algorithms

In this section, we present exploratory planning algorithms for low rank models, assuming that the dynamics are known. Formally, we consider an $H$ step low rank MDP $\widetilde{\mathcal{M}}$ with deterministic start state $x_0$, fixed action $a_0$, and transition matrices $T_0, \ldots, T_{H-1}$. Each transition operator $T_h$ factorizes as $T_h(x_{h+1} \mid x_h, a_h) = \langle \phi_h(x_h, a_h), \mu_h(x_{h+1})\rangle$ and we assume $\phi_{0:H-1}, \mu_{0:H-1}$ are *known*. To compartmentalize the results, we focus on exploratory planning at time $H$, but we will invoke these subroutines with MDP models that have horizon $h \leq H$. This simply requires rebinding variables.

We present two types of results. One style assumes that all expectations are computed exactly. As we are focusing purely on planning with known dynamics and rewards, this imposes a computational burden, but not a statistical one, while leading to a more transparent proof. To address the computational burden, we also consider algorithms that approximate all expectations with samples. For this, we assume that we can obtain sample transitions from the MDP model $\widetilde{\mathcal{M}}$ in a computationally efficient manner. Formally, the *sampling oracle* allows us to sample $x' \sim T_h(\cdot \mid x, a)$ for any $x, a$.

### C.1   Planning with a sampling oracle

For the computational style of result, it will be helpful to first show how to optimize a given reward function whenever the model admits a sampling oracle. As notation, we always consider an explicitly specified non-stationary reward function $R : \mathcal{X} \times \mathcal{A} \times \{0, \ldots, H-1\} \to [0,1]$. Then, we define

$$V(\pi, R) = \mathbb{E}\left[\sum_{h=0}^{H-1} R(x_h, a_h, h) \mid \pi, \widetilde{\mathcal{M}}\right].$$

The next lemma is a simple application of the result of Jin et al. (2019).

**Lemma 15.** *Suppose that the reward function $R : \mathcal{X} \times \mathcal{A} \times \{0, \ldots, H-1\} \to [0,1]$ is explicitly given and that $T_{0:H-1}$ is a known low rank MDP that enables efficient sampling. Then for any $\epsilon > 0$ there is an algorithm for finding a policy $\hat{\pi}$ such that with probability at least $1 - \delta$, $V(\hat{\pi}, R) \geq \max_\pi V(\pi, R) - \epsilon$ in polynomial time with $\mathrm{poly}(d, H, 1/\epsilon, \log(1/\delta))$ calls to the sampling routine.*

---

**Algorithm 3** Exploratory planner for simplex representations

---

**Input:** MDP $\widetilde{\mathcal{M}} = (\phi_{0:H-1}, \mu_{0:H-1})$ with $\phi_h(x_h, a_h) \in \Delta([d_{\mathrm{LV}}])$, $\mu_h[z] \in \Delta(\mathcal{X})$.

**for** $z = 1, \ldots, d_{\mathrm{LV}}$ **do**

  Compute $\pi_z = \mathrm{argmax}_\pi \mathbb{E}[\phi_{H-1}(x_{H-1}, a_{H-1})[z] \mid \widetilde{\mathcal{M}}, \pi]$

**end for**

Output policy mixture $\rho := \mathrm{unif}(\{\pi_z\}_{z=1}^{d_{\mathrm{LV}}})$

---

*Proof.* As we have sampling access to the MDP, we can execute the LSVI-UCB algorithm of Jin et al. (2019). For any $n$, if we execute the algorithm for $n$ episodes, it produces $n$ policies $\pi_1, \ldots, \pi_n$ and guarantees

$$\max_\pi V(\pi, R) - \frac{1}{n} \sum_{i=1}^n V(\pi_i, R) \leq c\sqrt{\frac{d^3 H^3 \log(ndH/\delta)}{n}}$$

with probability at least $1 - \delta$ where $c > 0$ is a universal constant. We are assured that one of the policies $\pi_1, \ldots, \pi_n$ is at most $\epsilon/2$-suboptimal by taking $n = O\left(d^3 H^3 \log(dH/(\epsilon\delta))/\epsilon^2\right)$.

We find this policy via a simple policy evaluation step. For each policy $\pi_i$, we collect $O(H^2 \log(n/\delta)/\epsilon^2)$ roll-outs using the generative model, where we take actions according to $\pi_i$. Via a union bound, this guarantees that for each $i$ we have $\hat{V}_i$ such that with probability at least $1 - \delta$

$$\max_i \left|\hat{V}_i - V(\pi_i, R)\right| \leq \epsilon/4.$$

Therefore, if we take $\hat{i} = \mathrm{argmax}_{i \in [n]} \hat{V}_i$ we are assured that $V(\pi_{\hat{i}}, R) \geq \max_\pi V(\pi, R) - \epsilon$ with probability at least $1 - 2\delta$. The total number of samples required from the model are

$$nH\left(1 + \frac{H^2 \log(n/\delta)}{\epsilon^2}\right) = \tilde{O}\left(\frac{d^3 H^6 \log(1/\delta)}{\epsilon^4}\right). \qquad \square$$

## C.2  Planning with simplex features

We first consider a simpler planning algorithm that is adapted to the simplex features representation. The pseudocode is displayed in Algorithm 3. The planner computes a mixture policy $\rho$, where component $\pi_i$ of the mixture focuses on activating coordinate $i$ of the feature map $\phi_{H-1}(x_{H-1}, a_{H-1})$. Each mixture component can be computed in a straightforward manner using a dynamic programming approach, such as LSVI. The basic guarantee for this algorithm is the following lemma.

**Lemma 16** (Guarantee for Algorithm 3). *If $\widetilde{\mathcal{M}}$ is an $H$-step low rank MDP with simplex features of dimension $d_{\mathrm{LV}}$, then the output of Algorithm 3, $\rho$, satisfies*

$$\forall \pi, z \in [d_{\mathrm{LV}}] : \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})[z] \mid \widetilde{\mathcal{M}}, \pi\right] \leq d_{\mathrm{LV}} \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})[z] \mid \widetilde{\mathcal{M}}, \rho\right].$$

*Given a sampling oracle for $\widetilde{\mathcal{M}}$, the algorithm runs in polynomial time with $\mathrm{poly}(d_{\mathrm{LV}}, H, 1/\varepsilon_{\mathrm{opt}}, \log(1/\delta))$ calls to SAMP, and with probability at least $1 - \delta$, $\rho$ satisfies*

$$\forall \pi, z \in [d_{\mathrm{LV}}], \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})[z] \mid \widetilde{\mathcal{M}}, \pi\right] \leq d_{\mathrm{LV}} \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})[z] \mid \widetilde{\mathcal{M}}, \rho\right] + \varepsilon_{\mathrm{opt}}.$$

*Proof.* The first result follows immediately from the non-negativity of $\phi_{H-1}(x_{H-1}, a_{H-1})[i]$, the optimality property of $\pi_i$ and the definition of $\rho$.

For the second result, by Lemma 15 we can optimize any explicitly specified reward function using a polynomial number of samples. If we call this sampling-based planner for each of the $d$ reward functions, with high probability (via a union bound) the policies $\hat{\pi}_i$ are near-optimal for their corresponding reward functions. By appropriately re-scaling the accuracy parameter in Lemma 15 we obtain the desired guarantee. $\qquad \square$

## C.3 Elliptical planner

The next planning algorithm applies to general low rank MDP, and it is more sophisticated. It proceeds in iterations, where in iteration $t$ we maintain a covariance matrix $\Sigma_{t-1}$ and, in (2), we search for a policy that maximizes quadratic forms with the inverse covariance $\Sigma_{t-1}^{-1}$. With a sampling oracle this optimization can be done via a call to Lemma 15. If this maximizing policy $\pi_t$ cannot achieve large quadratic forms against $\Sigma_{t-1}^{-1}$, then we halt and output the mixture of all previous policies. Otherwise, we mix $\pi_t$ into our candidate solution, update the covariance matrix accordingly, and advance to the next iteration. The performance guarantee for this algorithm is as follows.

**Lemma 17** (Guarantee for Algorithm 2). *If $\widetilde{\mathcal{M}}$ is an $H$-step low rank MDP with embedding dimension $d$ then for any $\beta > 0$, Algorithm 2 terminates after at most $T + 1$ iterations where $T \leq 4d \log(1 + 4/\beta)/\beta$. Upon termination, $\rho$ guarantees*

$$\forall \pi : \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})^\top (\Sigma_\rho + I/T)^{-1} \phi_{H-1}(x_{H-1}, a_{H-1}) \mid \widetilde{\mathcal{M}}, \pi\right] \leq T\beta.$$

*where $\Sigma_\rho = \frac{1}{T} \sum_{t=1}^T \Sigma_{\pi_t}$.*

*Proof.* The performance guarantee is immediate from the termination condition, using the fact that $\Sigma_T = T \cdot (\Sigma_\rho + I/T)$.

For the iteration complexity bound, we condense the notation and omit the dependence on $H - 1$, $x_{H-1}, a_{H-1}$ in all terms. We have

$$\beta T \leq \sum_{t=1}^T \mathbb{E}\left[\phi^\top \Sigma_{t-1}^{-1} \phi \mid \widetilde{\mathcal{M}}, \pi_t\right] = \sum_{t=1}^T \operatorname{tr}(\Sigma_{\pi_t} \Sigma_{t-1}^{-1}) \leq 2d \log(1 + T/d),$$

where the first inequality is based on the fact that we did not terminate at each iteration $t \in [T]$ and the last inequality follows from a standard elliptical potential argument (e.g., Lemma 11 in Dani et al. (2008); see Lemma 26 for a precise statement and proof). This gives an upper bound on $T$ that implies the one in the lemma statement, via Corollary 27. $\square$

With the sampling oracle, we modify the algorithm slightly and obtain a qualitatively similar guarantee. The modifications are discussed in the proof.

**Lemma 18.** *The sample-based version of Algorithm 2 has the following guarantee. Assume $\widetilde{\mathcal{M}}$ is an $H$-step low rank MDP with embedding dimension $d$ and fix $\beta > 0$, $\delta \in (0, 1)$. Then the algorithm terminates after at most $T + 1$ iterations, where $T \leq O(d \log(1 + 1/\beta)/\beta)$. Upon termination, it ouputs a matrix $\hat{\Sigma}$ and a policy $\rho$ such that with probability at least $1 - \delta$:*

$$\forall \pi : \mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})^\top \left(\hat{\Sigma} + I/T\right)^{-1} \phi_{H-1}(x_{H-1}, a_{H-1}) \mid \widetilde{\mathcal{M}}, \pi\right] \leq O(T\beta),$$

$$\left\|\hat{\Sigma} - \left(\mathbb{E}\left[\phi_{H-1}(x_{H-1}, a_{H-1})\phi_{H-1}(x_{H-1}, a_{H-1})^\top \mid \rho, \widetilde{\mathcal{M}}\right] + I/T\right)\right\|_{\operatorname{op}} \leq O(\beta/d).$$

*The algorithm runs in polynomial time with $\operatorname{poly}(d, H, 1/\beta, \log(1/\delta))$ calls to the sampling oracle.*

*Proof.* The algorithm is modified as follows. We replace all covariances with empirical approximations, obtained by calls to the sampling subroutine. We call the empirical versions $\hat{\Sigma}_t, \hat{\Sigma}_{\pi_t}$, etc. Then, the policy optimization step (2) is performed via an application of Lemma 15 and so we find an $\varepsilon_{\operatorname{opt}}$-suboptimal policy $\pi_t$ for the reward function induced by $\hat{\Sigma}_{t-1}$. Then we use the sampling subroutine to estimate the value of this policy, which we denote $\hat{V}_t(\pi_t)$. As before, we terminate if $\hat{V}_t(\pi_t) \leq \beta$. If we terminate in round $t$, we

output $\rho = \mathrm{unif}(\{\pi_i\}_{i=1}^{t-1})$ and we also output $\hat{\Sigma} = \frac{1}{t-1}\sum_{i=1}^{t-1}\hat{\Sigma}_{\pi_i}$. As notation, we use $V_t(\pi)$ to denote the value for policy $\pi$ on the reward function used in iteration $t$, which is induced by $\hat{\Sigma}_{t-1}$.

With $\mathrm{poly}(d, H, T, 1/\varepsilon_{\mathrm{opt}}, \log(1/\delta))$ calls to the sampling subroutine and assuming the total number of iterations of the algorithm $T$ is polynomial, we can verify that with probability $1 - \delta$

$$\max_{t \in [T]} \max \left\{ d \cdot \left\| \hat{\Sigma}_{\pi_t} - \Sigma_{\pi_t} \right\|_{\mathrm{op}}, \left| \hat{V}_t(\pi_t) - V_t(\pi_t) \right|, \max_\pi V_t(\pi) - V_t(\pi_t) \right\} \leq \varepsilon_{\mathrm{opt}}.$$

The first two bounds follow from standard concentration of measure arguments. The final one is based on an application of Lemma 15.

Now, if we terminate in iteration $t$, we know that $\hat{V}_t(\pi_t) \leq \beta$. This implies

$$\max_\pi V_t(\pi) \leq V_t(\pi_t) + \varepsilon_{\mathrm{opt}} \leq \hat{V}_t(\pi_t) + 2\varepsilon_{\mathrm{opt}} \leq \beta + 2\varepsilon_{\mathrm{opt}}.$$

As we are interested in the reward function induced by $\hat{\Sigma}_{t-1}$, this verifies the quality guarantee, provided $\varepsilon_{\mathrm{opt}} = O(\beta)$.

Finally, we turn to the iteration complexity. Similarly to above, we have

$$T(\beta - 2\varepsilon_{\mathrm{opt}}) \leq \sum_{t=1}^T \hat{V}_t(\pi_t) - 2\varepsilon_{\mathrm{opt}} \leq \sum_{t=1}^T V_t(\pi_t) - \varepsilon_{\mathrm{opt}}$$

$$= \sum_{t=1}^T \mathbb{E}\left[ \phi^\top \hat{\Sigma}_{t-1}^{-1} \phi \mid \widetilde{\mathcal{M}}, \pi_t \right] - \varepsilon_{\mathrm{opt}} = \sum_{t=1}^T \mathrm{tr}(\Sigma_{\pi_t} \hat{\Sigma}_{t-1}^{-1}) - \varepsilon_{\mathrm{opt}}$$

$$\leq \sum_{t=1}^T \mathrm{tr}(\hat{\Sigma}_{\pi_t} \hat{\Sigma}_{t-1}^{-1}) \leq 2d\log(1 + T/d).$$

In other words, if we set $\varepsilon_{\mathrm{opt}} = O(\beta)$ then both the iteration complexity and the performance guarantee are unchanged. The accuracy guarantee for the covariance matrix $\hat{\Sigma}_{t-1}$ is straightforward, since each $\hat{\Sigma}_{\pi_t}$ is $\varepsilon_{\mathrm{opt}}$ accurate and $\hat{\Sigma}$ is the average of such matrices. $\qquad\square$

# D  Planning in the environment

In this section, we prove Theorem 4, which is based on planning in the environment, rather than in the model. Recall that the advantage of this approach is that we do not need the sampling oracle, SAMP, but the downside is that we require Assumption 2. As Assumption 2 implies a polynomial upper bound on $d_{\mathrm{LV}}$, we focus on the simplex representation, and remark that this also accommodates general representations with polynomial overhead in sample complexity. Planning in the environment with general representations (and Assumption 2) is possible, but the arguments and calculations are much simpler in the simplex case.

The algorithm here is a "forward" version of FLAMBE, that we call FLAMBE.F. The pseudocode is displayed in Algorithm 4. The algorithm learns the dynamics one time step at a time, starting from $h = 0$ to $h = H - 1$. In iteration $h$, we use an exploratory policy $\rho_h$ to collect a dataset of triples $(x_h, a_h, x_{h+1})$ that we pass to MLE to obtain an estimate $\hat{T}_h$ of the dynamics at time $h$. Then we pass the previously computed policies and feature maps to a new planning algorithm (Algorithm 5), which yields the exploratory policy $\rho_{h+1}$ for the next iteration.

The main argument is based on inductively establishing two facts.

$$\forall h' < h, \forall \pi : \ \mathbb{E}\left[ \left\| \hat{T}_{h'}(\cdot \mid x_{h'}, a_{h'}) - T_{h'}(\cdot \mid x_{h'}, a_{h'}) \right\|_{\mathrm{TV}}^2 \mid \pi, \mathcal{M} \right] \leq \varepsilon_{\mathrm{TV}} \tag{6}$$

$$\forall z \in \mathcal{Z}_h : \ \max_\pi \mathbb{P}\left[ z_h = z \mid \pi, \mathcal{M} \right] \leq \kappa \cdot \mathbb{P}\left[ z_h = z \mid \rho_h, \mathcal{M} \right], \tag{7}$$

---

**Algorithm 4** FLAMBE.F: Forward version of FLAMBE with simplex features

---

**Input:** Environment $\mathcal{M}$, function classes $\Phi, \Upsilon$, subroutine MLE, parameter $n$.

Set $\rho_0$ to be the null policy, which takes no actions.

**for** $h = 0, \ldots, H-1$ **do**

  Set $\rho_h^{\text{train}} \leftarrow \rho_h \circ \text{unif}(\mathcal{A})$.         {Uniform over available actions.}

  Collect $n$ triples $D_h \leftarrow \{(x_h^{(i)}, a_h^{(i)}, x_{h+1}^{(i)})\}_{i=1}^n$ by executing $\rho_h^{\text{train}}$ in $\mathcal{M}$.

  Solve maximum likelihood problem: $(\hat{\phi}_h, \hat{\mu}_h) \leftarrow \text{MLE}(D_h)$.

  Set $\hat{T}_h(x_{h+1} \mid x_h, a_h) = \left\langle \hat{\phi}_h(x_h, a_h), \hat{\mu}_h(x_{h+1}) \right\rangle$.

  Call planner (Algorithm 5) with policies $\rho_{0:h}$ and feature maps $\hat{\phi}_{0:h-1}$ to obtain $\rho_h^{\text{pre}}$.

  Set $\rho_{h+1} = \rho_h^{\text{pre}} \circ \text{unif}(\mathcal{A})$.

**end for**

---

where $\kappa > 0$ is a constant we will set towards the end of the proof. We do this in two parts. We first focus on optimizing a fixed given reward function by collecting experience from the environment, analogously to the sampling approach in Lemma 15. For this part, we will assume that (6) and (7) hold with $h$ being the current planning horizon. In the next subsection we choose the reward functions carefully to establish the guarantees required by FLAMBE. Since we are considering simplex representations, this second part is very similar to Algorithm 3.

## D.1 Optimizing a fixed reward function

To optimize a fixed reward function, the high level idea is that, via Lemma 1, we can approximate any Bellman backup using our features $\hat{\phi}$, and via (7), we can collect a dataset with good coverage. Using these two properties the planning algorithm, LINEAR-FQI, displayed as a subroutine in Algorithm 5 is quite natural. The algorithm is a least squares dynamic programming algorithm (FQI stands for "Fitted Q Iteration"). For each $h$, working from $H-1$ down to 0, we collect a dataset of $n$ samples by following $\rho_h$. Then, we solve a least squares regression problem to approximate the Bellman backup of the value function estimate $\hat{V}_{h+1}$ for the next time. We use this to define the value function and the policy for the current time in the obvious way.

  Note that we index policies in two different ways: $\rho_h$ is the exploratory policy that induces a distribution over $x_h$, while $\hat{\pi}_h$ is the one-step policy that we acts on $x_h$. As with the other planning lemmas, we apply the next lemma with a value of $H$ that is not necessarily the real horizon in the environment. In particular, we will use this lemma in the $h^{\text{th}}$ iteration of FLAMBE, with planning horizon $h-1$ and with reward functions specified in the next subsection. By induction, we can assume that (6) and (7) hold for the planning horizon.

**Lemma 19.** *Assume that* (6) *and* (7) *hold for all* $h \in [H]$. *Then for any reward functions* $R_{0:H-1} : \mathcal{X} \times \mathcal{A} \rightarrow [0,1]$ *and any* $\delta \in (0,1)$, *if we set*

$$n \geq \frac{2304d^3}{\varepsilon_{\text{TV}}^2} \log\left(1152d^3/\varepsilon_{\text{TV}}^2\right) + \frac{2304d^2}{\varepsilon_{\text{TV}}^2} \log(2H/\delta),$$

*then the policy* $\hat{\pi}_{0:H-1}$ *returned by* LINEAR-FQI *satisfies*

$$\mathbb{E}\left[\sum_{h=0}^{H-1} r_h \mid \hat{\pi}_{0:H-1}, \mathcal{M}\right] \geq \max_{\pi} \mathbb{E}\left[\sum_{h=0}^{H-1} r_h \mid \pi, \mathcal{M}\right] - 2H^3\sqrt{2\kappa K \varepsilon_{\text{TV}}}.$$

*Proof.* The analysis is similar to that of Chen and Jiang (2019), who study a similar algorithm in the infinite-horizon discounted setting. Let $\mathbb{E}_h$ denote expectation induced by the distribution over $(x_h, a_h, x_{h+1})$

---

**Algorithm 5** Planning in the environment with simplex features

---

**input:** Exploratory policies $\rho_{0:H-1}$, feature maps $\hat{\phi}_{0:H-1}$ with $\hat{\phi}_h(x,a) \in \Delta([d_{\mathrm{LV}}])$.
**for** $i = 1, \ldots, d_{\mathrm{LV}}$ **do**
      Compute $\hat{\pi}_i = \text{LINEAR-FQI}(n, \rho_{0:H-1}, \hat{\phi}_{0:H-1}, R_{H-1} := \hat{\phi}_{H-1}(x,a)[i])$      $\{R_{0:H-2} \equiv 0\}$
**end for**
**return** policy mixture $\rho := \text{unif}(\{\hat{\pi}_i\}_{i=1}^{d_{\mathrm{LV}}})$.


**function** LINEAR-FQI$(n, \rho_{0:H-1}, \hat{\phi}_{0:H-1}, R_{0:H-1})$
**input:** Sample size $n$, policies $\rho_{0:H-1}$, feature maps $\hat{\phi}_{0:H-1}$, rewards $R_{0:H-1} : \mathcal{X} \times \mathcal{A} \to [0,1]$.
Set $\hat{V}_H(x) = 0$
**for** $h = H-1, \ldots, 0$ **do**
      Collect $n$ samples $\{(x_h^{(i)}, a_h^{(i)}, x_{h+1}^{(i)})\}_{i=1}^n$ by following $\rho_h \circ \text{unif}(\mathcal{A})$ in $\mathcal{M}$.
      Solve least squares problem:

$$\hat{\theta}_h \leftarrow \underset{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq H\sqrt{d}}{\operatorname{argmin}} \sum_{i=1}^n \left( \left\langle \theta, \hat{\phi}_h(x_h^{(i)}, a_h^{(i)}) \right\rangle - \hat{V}_{h+1}(x_{h+1}^{(i)}) \right)^2.$$

      Define $\hat{Q}_h(x,a) = R_h(x,a) + \left\langle \hat{\theta}_h, \hat{\phi}_h(x,a) \right\rangle$.
      Define $\hat{\pi}_h(x) = \operatorname{argmax}_a \hat{Q}_h(x,a), \hat{V}_h(x) = \min\{\max_a \hat{Q}_h(x,a), H\}$.
**end for**
**return** $\hat{\pi} = (\hat{\pi}_0, \ldots, \hat{\pi}_{H-1})$.

---

obtained by following $\rho_h \circ \text{unif}(\mathcal{A})$. For any function $f : \mathcal{X} \to \mathbb{R}$, let $\mathcal{B}_h f(x,a) := \mathbb{E}[f(x_{h+1}) \mid x_h, a_h]$ denote the Bellman backup operator for time $h$ *without the immediate reward*. Let $\hat{\mathcal{B}}_h$ denote the Bellman backup operator induced by the learned model at time $h$, again without the immediate reward. We omit the dependence on $x, a$ in these operators when it is clear from context. Note that by the normalization assumptions, we always have $\hat{V}_{h+1}(x_{h+1}) \in [0, H]$. Moreover, $\hat{\mathcal{B}}_h \hat{V}_{h+1}$ is a linear function in $\hat{\phi}_h$ where the coefficient vector has $\ell_2$ norm at most $H\sqrt{d}$.

We apply Lemma 20 with $B := H\sqrt{d}$, and we take a union bound over all $h \in [H]$. Defining $\Delta_n := 24H^2 d\sqrt{2(d \log n + \log(2H/\delta))/n}$, we have that with probability at least $1 - \delta$, for all $h \in [H]$:

$$\mathbb{E}_h \left[ \left( \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle - \mathcal{B}_h \hat{V}_{h+1} \right)^2 \right] \leq \min_{\theta : \|\theta\|_2 \leq B} \mathbb{E}_h \left[ \left( \left\langle \theta, \hat{\phi}_h(x_h, a_h) \right\rangle - \mathcal{B}_h \hat{V}_{h+1} \right)^2 \right] + \Delta_n$$

$$\leq \mathbb{E}_h \left[ \left( \hat{\mathcal{B}}_h \hat{V}_{h+1} - \mathcal{B}_h \hat{V}_{h+1} \right)^2 \right] + \Delta_n$$

$$\leq H^2 \cdot \mathbb{E}_h \left\| \hat{T}_h(\cdot \mid x_h, a_h) - T_h(\cdot \mid x_h, a_h) \right\|_{\mathrm{TV}}^2 + \Delta_n.$$

The first inequality here is the least squares generalization analysis, additionally using that $\mathcal{B}_h \hat{V}_{h+1}$ is the Bayes optimal predictor. The second uses the fact that the Bellman backups in the model are linear functions in $\hat{\phi}$ (with bounded coefficient vector). Precisely, we have

$$\hat{\mathcal{B}}_h \hat{V}_{h+1}(x_h, a_h) = \left\langle \hat{\phi}_h(x_h, a_h), \int \hat{\mu}_h(x_{h+1}) \hat{V}_{h+1}(x_{h+1}) d(x_{h+1}) \right\rangle.$$

Setting $\theta$ to be the second term, we obtain the second inequality. Finally, we apply Holder's inequality and

use the fact that $\hat{V}_{h+1}$ is bounded in $[0, H]$ by construction. Appealing to (6) we have

$$\mathbb{E}_h \left[ \left( \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle - \mathcal{B}_h \hat{V}_{h+1} \right)^2 \right] \leq H^2 \varepsilon_{\mathrm{TV}} + \Delta_n.$$

Now applying (7), we transfer this squared error to the distribution induced by any other policy. This calculation is exactly the same as in the main induction argument (see (9)), and it yields

$$\mathbb{E}_\pi \left[ \left( \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle - \mathcal{B}_h \hat{V}_{h+1} \right)^2 \right] \leq \kappa K \left( H^2 \varepsilon_{\mathrm{TV}} + \Delta_n \right).$$

Next, we bound the difference in cumulative rewards between $\hat{\pi} := \hat{\pi}_{0:H-1}$ and the optimal policy $\pi^\star$ for the reward function. For this, recall that we define $\hat{Q}_0(x, a) = R_0(x, a) + \langle \hat{\theta}_0, \hat{\phi}_0(x, a) \rangle$ and also that $\hat{\pi}_0$ is greedy with respect to this $Q$ function, which implies that $\hat{Q}_0(x, \hat{\pi}_0(x)) \geq \hat{Q}_0(x, \pi^\star(x))$ for all $x$. Therefore,

$$
\begin{aligned}
V^\star - V^{\hat{\pi}} &= \mathbb{E} \left[ R(x_0, a_0) + V^\star(x_1) \mid \pi^\star \right] - \mathbb{E} \left[ R(x_0, a_0) + V^{\hat{\pi}}(x_1) \mid \hat{\pi} \right] \\
&\leq \mathbb{E} \left[ R(x_0, a_0) + V^\star(x_1) - \hat{Q}_0(x_0, a_0) \mid \pi^\star \right] - \mathbb{E} \left[ R(x_0, a_0) + V^{\hat{\pi}}(x_1) - \hat{Q}_0(x_0, a_0) \mid \hat{\pi} \right] \\
&= \mathbb{E} \left[ V^\star(x_1) - \left\langle \hat{\theta}_0, \hat{\phi}_0(x_0, a_0) \right\rangle \mid \pi^\star \right] - \mathbb{E} \left[ V^{\hat{\pi}}(x_1) - \left\langle \hat{\theta}_0, \hat{\phi}_0(x_0, a_0) \right\rangle \mid \hat{\pi} \right] \\
&= \mathbb{E} \left[ V^\star(x_1) - \left\langle \hat{\theta}_0, \hat{\phi}_0(x_0, a_0) \right\rangle \mid \pi^\star \right] - \mathbb{E} \left[ V^\star(x_1) - \left\langle \hat{\theta}_0, \hat{\phi}_0(x_0, a_0) \right\rangle \mid \hat{\pi} \right] \\
&\quad + \mathbb{E} \left[ V^\star(x_1) - V^{\hat{\pi}}(x_1) \mid \hat{\pi} \right].
\end{aligned}
$$

Continuing, we find

$$
\begin{aligned}
V^\star - V^{\hat{\pi}} &\leq \sum_{h=0}^{H-1} \mathbb{E} \left[ V^\star(x_{h+1}) - \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle \mid \hat{\pi}_{0:h-1} \circ \pi^\star \right] \\
&\quad - \sum_{h=0}^{H-1} \mathbb{E} \left[ V^\star(x_{h+1}) - \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle \mid \hat{\pi}_{0:h} \right].
\end{aligned}
$$

Next, we bound each of these terms. Let us focus on just one of them, call the roll-in policy $\pi$ and drop the dependence on $h$. Then,

$$
\begin{aligned}
\mathbb{E}_\pi \left[ \left| \mathbb{E} \left[ V^\star(x') \mid x, a \right] - \left\langle \hat{\theta}, \hat{\phi}(x, a) \right\rangle \right| \right] &\leq \mathbb{E}_\pi \left[ \left| \mathcal{B} V^\star(x, a) - \mathcal{B} \hat{V}(x, a) \right| + \left| \mathcal{B} \hat{V}(x, a) - \left\langle \hat{\theta}, \hat{\phi}(x, a) \right\rangle \right| \right] \\
&\leq \mathbb{E}_\pi \left[ \left| V^\star(x') - \hat{V}(x') \right| + \left| \mathcal{B} \hat{V}(x, a) - \left\langle \hat{\theta}, \hat{\phi}(x, a) \right\rangle \right| \right],
\end{aligned}
$$

where the second inequality is Jensen's inequality. By definition

$$
\begin{aligned}
\mathbb{E}_\pi \left[ \left| V^\star(x') - \hat{V}(x') \right| \right] &= \mathbb{E}_\pi \left[ \left| \max_a Q^\star(x', a) - \min\{H, \max_{a'} R(x', a') + \left\langle \hat{\phi}(x', a'), \hat{\theta} \right\rangle \} \right| \right] \\
&\leq \mathbb{E}_\pi \left[ \left| \max_a Q^\star(x', a) - \max_{a'} R(x', a') + \left\langle \hat{\phi}(x', a'), \hat{\theta} \right\rangle \right| \right] \leq \mathbb{E}_{\pi \circ \tilde{\pi}} \left[ \left| \mathcal{B} V^\star(x', a') - \left\langle \hat{\phi}(x', a'), \hat{\theta} \right\rangle \right| \right].
\end{aligned}
$$

Here in the last inequality, we define $\tilde{\pi}$ to choose the larger of the two actions, that is we set $\tilde{\pi}(x') = \mathrm{argmax}_{a \in \mathcal{A}} \max\{Q^\star(x', a), R(x', a) + \langle \hat{\theta}, \hat{\phi}(x', a) \rangle\}$. This expression has the same form as the initial one, but at the next time point, so unrolling, we get

$$
\begin{aligned}
\mathbb{E} \left[ \left| \mathcal{B} V^\star(x_h, a_h) - \left\langle \hat{\theta}_h, \hat{\phi}_h(x_h, a_h) \right\rangle \right| \mid \pi \right] &\leq \sum_{\tau=h}^{H-1} \max_{\pi_\tau} \mathbb{E}_{\pi_\tau} \left[ \left| \mathcal{B} \hat{V}_{\tau+1}(x_\tau, a_\tau) - \left\langle \hat{\theta}_\tau, \hat{\phi}_\tau(x_\tau, a_\tau) \right\rangle \right| \right] \\
&\leq H \sqrt{\kappa K (H^2 \varepsilon_{\mathrm{TV}} + \Delta_n)}.
\end{aligned}
$$

Plugging this into the overall value difference, the final bound is

$$V^\star - V^{\hat{\pi}} \leq 2H^2 \sqrt{\kappa K (H^2 \varepsilon_{\mathrm{TV}} + \Delta_n)}.$$

To wrap up, we want the term involving $\Delta_n$ to be at most $H^2 \varepsilon_{\mathrm{TV}}$, so the term involving $\Delta_n$ is of the same order as the term involving $\varepsilon_{\mathrm{TV}}$. By our definition of $\Delta_n$, this requires

$$n \geq \frac{24^2 d^2}{\varepsilon_{\mathrm{TV}}^2} \cdot 2 \left( d \log n + \log(2H/\delta) \right).$$

A sufficient condition here is

$$n \geq \frac{2304 d^3}{\varepsilon_{\mathrm{TV}}^2} \log \left( 1152 d^3 / \varepsilon_{\mathrm{TV}}^2 \right) + \frac{2304 d^2}{\varepsilon_{\mathrm{TV}}^2} \log(2H/\delta),$$

which yields the result. $\qquad\square$

**Lemma 20.** *Let $\{\phi_i, y_i\}_{i=1}^n$ be $n$ samples drawn iid from some distribution where $\phi \in \mathbb{R}^d$ satisfies $\|\phi\|_2 \leq 1$ and $y \in [0, H]$ almost surely. Let $\hat{\theta} \in \mathbb{R}^d$ denote the constrained square loss minimizer, constrained so that $\|\hat{\theta}\|_2 \leq B$, where $B \geq H$. Then for any $\delta \in (0,1)$, with probability at least $1 - \delta$ we have*

$$\mathbb{E}\left[ (\langle \hat{\theta}, \phi \rangle - y)^2 \right] \leq \min_{\theta : \|\theta\|_2 \leq B} \mathbb{E}\left[ (\langle \theta, \phi \rangle - y)^2 \right] + 24 B^2 \sqrt{\frac{2}{n} \left( d \log(n) + \log(2/\delta) \right)}.$$

*Proof.* Fix $\theta$ with $\|\theta\|_2 \leq B$. We will apply Hoeffding's inequality on this $\theta$ and then use a covering argument for uniform convergence. Let $R(\theta)$ denote the expected square loss, with $\hat{R}(\theta)$ as the empirical counterpart. Using the bounds on all quantities, the square loss has range $(B + H)^2 \leq 4B^2$, and so Hoeffding's inequality yields that with probability at least $1 - \delta$

$$\left| R(\theta) - \hat{R}(\theta) \right| \leq 4B^2 \sqrt{\frac{2}{n} \log(2/\delta)}.$$

Let $V_\gamma$ denote a covering of $\{\theta : \|\theta\|_2 \leq B\}$ in the $\ell_2$ norm at scale $\gamma$, which has $\log |V_\gamma| \leq d \log(2B/\gamma)$ via standard arguments. Taking a union bound, the above inequality holds for all $\theta \in V_\gamma$ with probability $1 - |V_\gamma| \delta$. By direct calculation, we see that $\hat{R}(\theta)$ and $R(\theta)$ are both $2(B + H)$-Lipschitz. Therefore, we have that with probability $1 - |V_\gamma| \delta$, for all $\theta$ with $\|\theta\|_2 \leq B$

$$\left| R(\theta) - \hat{R}(\theta) \right| \leq 4B\gamma + 4B^2 \sqrt{\frac{2}{n} \log(2/\delta)}.$$

Taking $\gamma = 2B/\sqrt{n}$ we can rebind $\delta$ and absorb the first term into the second. Thus, with probability at least $1 - \delta$, for all $\theta$, we have

$$\left| R(\theta) - \hat{R}(\theta) \right| \leq 12 B^2 \sqrt{\frac{2}{n} \left( d \log(n) + \log(2/\delta) \right)}.$$

Now by the standard ERM analysis, we have

$$R(\hat{\theta}) \leq \hat{R}(\hat{\theta}) + 12 B^2 \sqrt{\frac{2}{n} \left( d \log(n) + \log(2/\delta) \right)} \leq \min_\theta \hat{R}(\theta) + 12 B^2 \sqrt{\frac{2}{n} \left( d \log(n) + \log(2/\delta) \right)}$$

$$\leq \min_\theta R(\theta) + 24 B^2 \sqrt{\frac{2}{n} \left( d \log(n) + \log(2/\delta) \right)}. \qquad\square$$

## D.2 Instantiating the reward functions.

We now use Algorithm 5 in FLAMBE.F, to establish the induction. Assume that $\hat{\phi}_h(x,a) \in \Delta([d_{\mathrm{LV}}])$ for all $x, a, h$, analogously to in Theorem 3. Recall that $\rho_h$ is our exploratory policy that induces a distribution over $x_h$. We augment $\rho_h$ with an action taken uniformly at random to obtain the "training policy" $\rho_h^{\mathrm{train}}$. Via an application of Theorem 21 (using Assumption 1), we know that with probability at least $1 - \delta$ we learn $\hat{\phi}_h, \hat{\mu}_h$ such that (with $\hat{T} = \langle \hat{\phi}_h, \hat{\mu}_h \rangle$):

$$\mathbb{E}_{(x_h,a_h)\sim\rho_h^{\mathrm{train}}} \left\| \hat{T}_h(\cdot \mid x_h, a_h) - T_h(\cdot \mid x_h, a_h) \right\|_{\mathrm{TV}}^2 \leq \frac{2\log(|\Phi||\Upsilon|/\delta)}{n} =: \varepsilon_{\mathrm{sup}}. \tag{8}$$

This is the only step where we use the optimization oracle, MLE, and similar guarantee can also be obtained by other means. As one example, in Remark 22, we discuss a generative adversarial approach.

We now use this bound and (7) to establish (6) for time $h$. Considering any policy $\pi$, we define the "error function" $\mathrm{err}_\pi(x_h) := \int \pi(a_h \mid x_h) \left\| \hat{T}_h(\cdot \mid x_h, a_h) - T_h(\cdot \mid x_h, a_h) \right\|_{\mathrm{TV}}^2$.

$$
\begin{aligned}
\mathbb{E}_\pi &\left\| \hat{T}_h(\cdot \mid x_h, a_h) - T_h(\cdot \mid x_h, a_h) \right\|_{\mathrm{TV}}^2 = \mathbb{E}_\pi \left[ \mathrm{err}_\pi(x_h) \right] \\
&= \sum_{z \in \mathcal{Z}_h} \mathbb{P}\left[ z_h = z \mid \pi, \mathcal{M} \right] \cdot \int \mathrm{err}_\pi(x_h) \nu_{h-1}^\star(x_h \mid z) d(x_h) \\
&\leq \kappa \cdot \sum_{z \in \mathcal{Z}_h} \mathbb{P}\left[ z_h = z \mid \rho_h, \mathcal{M} \right] \cdot \int \mathrm{err}_\pi(x_h) \nu_{h-1}^\star(x_h \mid z) d(x_h) \\
&= \kappa \cdot \mathbb{E}_{\rho_h} \left[ \mathrm{err}_\pi(x_h) \right] \leq \kappa \cdot \mathbb{E}_{\rho_h^{\mathrm{train}}} \left[ \left\| \hat{T}_h(\cdot \mid x_h, a_h) - T_h(\cdot \mid x_h, a_h) \right\|_{\mathrm{TV}}^2 \right] \cdot \sup_{x_h, a_h} \left| \frac{\pi(a_h \mid x_h)}{\rho_h^{\mathrm{train}}(a_h \mid x_h)} \right| \\
&\leq \kappa K \varepsilon_{\mathrm{sup}} =: \varepsilon_{\mathrm{TV}}
\end{aligned}
\tag{9}
$$

The first inequality is (7), which allows us to transfer from the distribution induced by $\pi$ to the distribution induced by $\rho_h$. It is crucial that the pre-multiplier term involving $\nu_{h-1}^\star$ and $\mathrm{err}_\pi$ is non-negative which follows from the fact that $\mathrm{err}_\pi$ is non-negative and $\nu_{h-1}^\star(\cdot)[i]$ is a (positive) measure. The final two inequalities are based on importance weighting for the action at time $h$, using the fact that $\rho_h^{\mathrm{train}}(\cdot \mid x_h) = \mathrm{unif}(\mathcal{A})$. This final expression is our choice of $\varepsilon_{\mathrm{TV}}$, which establishes (6) for time $h$. For time $h = 0$, (6) follows immediately from (8), since $(x_0, a_0)$ are fixed. In particular all policies induce the same distribution over $(x_0, a_0)$ so transfering from $\pi$ to $\rho_0^{\mathrm{train}}$ is trivial. As $K, \kappa \geq 1$, this gives the base case.

Next we turn to establishing (7) for the next exploratory policy $\rho_{h+1}$. The planning algorithm in Algorithm 5 is analogous to Algorithm 3, except that we perform the optimization in the environment using LINEAR-FQI, with parameter $n$ that we will set subsequently. At iteration $h$ of FLAMBE.F, this yields $d_{\mathrm{LV}}$ policies $\hat{\pi}_1, \ldots, \hat{\pi}_{d_{\mathrm{LV}}}$ where $\hat{\pi}_i$ approximately maximizes the probability of reaching the $i^{\mathrm{th}}$ coordinate of $\hat{\phi}_{h-1}$ when executed in the real world.

Defining $\varepsilon_{\mathrm{stat}}$ to be the sub-optimality guaranteed by Lemma 19 (additionally taking a union bound over all $Hd_{\mathrm{LV}}$ invocations), we have that at iteration $h$ of FLAMBE

$$\mathbb{E}\left[ \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[i] \mid \hat{\pi}_i, \mathcal{M} \right] \geq \max_\pi \mathbb{E}\left[ \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[i] \mid \pi, \mathcal{M} \right] - \varepsilon_{\mathrm{stat}}.$$

We define $\rho_h^{\mathrm{pre}}$ to be the uniform distribution over the $\hat{\pi}_i$ policies, which induce a distribution over $x_h$.

Now for any function $f : \mathcal{X} \to [0,1]$, appealing to (6) at time $h$ we have

$$
\mathbb{E}_\pi f(x_h) \leq \mathbb{E}_\pi \left\langle \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}), \int \hat{\mu}_{h-1}(x_h) f(x_h) \right\rangle + \sqrt{\varepsilon_{\mathrm{TV}}}
$$

$$
= \sum_{i=1}^{d_{\mathrm{LV}}} \left( \int \hat{\mu}_{h-1}(x_h)[i] f(x_h) \right) \cdot \mathbb{E}_\pi \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[i] + \sqrt{\varepsilon_{\mathrm{TV}}}
$$

$$
\leq \sum_{i=1}^{d_{\mathrm{LV}}} \left( \int \hat{\mu}_{h-1}(x_h)[i] f(x_h) \right) \cdot \left( \mathbb{E}_{\hat{\pi}_i} \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[i] + \varepsilon_{\mathrm{stat}} \right) + \sqrt{\varepsilon_{\mathrm{TV}}}
$$

$$
\leq \sum_{i=1}^{d_{\mathrm{LV}}} \left( \int \hat{\mu}_{h-1}(x_h)[i] f(x_h) \right) \cdot \left( d_{\mathrm{LV}} \mathbb{E}_{\rho_h^{\mathrm{pre}}} \hat{\phi}_{h-1}(x_{h-1}, a_{h-1})[i] + \varepsilon_{\mathrm{stat}} \right) + \sqrt{\varepsilon_{\mathrm{TV}}}
$$

$$
\leq d_{\mathrm{LV}} \mathbb{E}_{\rho_h^{\mathrm{pre}}} \left\langle \hat{\phi}_{h-1}(x_{h-1}, a_{h-1}), \int \hat{\mu}_{h-1}(x_h)[i] f(x_h) \right\rangle + d_{\mathrm{LV}} \varepsilon_{\mathrm{stat}} + \sqrt{\varepsilon_{\mathrm{TV}}}
$$

$$
\leq d_{\mathrm{LV}} \mathbb{E}_{\rho_h^{\mathrm{pre}}} f(x_h) + d_{\mathrm{LV}} \varepsilon_{\mathrm{stat}} + (1 + d_{\mathrm{LV}}) \sqrt{\varepsilon_{\mathrm{TV}}}.
$$

The first and last inequalities here use (6) on $\hat{T}_{h-1}$, which holds by induction. The second inequality is the optimality guarantee for $\hat{\pi}_i$, and the third is based on the definition of $\rho_h^{\mathrm{pre}}$. For the fourth inequality, we collect terms, and additionally use that $f$ is $\ell_\infty$ bounded and $\hat{\mu}_{h-1}[i]$ is a measure, so $\left| \int \hat{\mu}_{h-1}(x_h)[i] f(x_h) \right| \leq 1$. Via importance weighting, we have that for any latent variable $z \in \mathcal{Z}_{h+1}$

$$
\max_\pi \mathbb{P}[z_{h+1} = z] \leq d_{\mathrm{LV}} K \cdot \mathbb{P}_{\rho_{h+1}}[z_{h+1} = z] + d_{\mathrm{LV}} \varepsilon_{\mathrm{stat}} + (1 + d_{\mathrm{LV}}) \sqrt{\varepsilon_{\mathrm{TV}}}.
$$

We must set the additive error to be at most $\eta_{\min}/2$, which establish (7) with $\kappa = 2 d_{\mathrm{LV}} K$. Unpacking the definition of $\varepsilon_{\mathrm{stat}}$ in the simplex features case this gives the constraint

$$
2H^3 \sqrt{4 d_{\mathrm{LV}} K^2 \varepsilon_{\mathrm{TV}}} + (1 + d_{\mathrm{LV}}) \sqrt{\varepsilon_{\mathrm{TV}}} \leq \eta_{\min}/2.
$$

Therefore, we set

$$
\varepsilon_{\mathrm{TV}} \leq O \left( \frac{\eta_{\min}^2}{H^6 d_{\mathrm{LV}}^2 K^2} \right).
$$

With this choice we establish the inductive hypothesis for the next time, and proceeding in this manner, we establish (6) for all time steps. As this is quite similar to our desired system identification guarantee, (1), we are just left to set the parameters appropriately and calculate the sample complexity. For the calls to MLE, we have that the $\varepsilon_{\mathrm{TV}} = 2 d_{\mathrm{LV}} K^2 \varepsilon_{\mathrm{sup}}$ which yields the constraint

$$
\varepsilon_{\mathrm{sup}} \leq O \left( \min \left\{ \frac{\eta_{\min}^2}{H^6 d_{\mathrm{LV}}^3 K^4}, \frac{\varepsilon^2}{d_{\mathrm{LV}} K^2} \right\} \right).
$$

This means that for the calls to MLE we may set $n$ as

$$
n = O \left( \max \left\{ \frac{d_{\mathrm{LV}}^2 K^2 H^6}{\eta_{\min}^2}, \frac{1}{\varepsilon^2} \right\} \cdot d_{\mathrm{LV}} K^2 \log(H |\Phi| |\Upsilon| / \delta) \right).
$$

The calls to MLE incur a total sample complexity of $nH$.

We also have to collect trajectories to invoke LINEAR-FQI. For this, we must set $n$ as

$$
n = \tilde{O} \left( \frac{d_{\mathrm{LV}}^3}{\varepsilon_{\mathrm{TV}}^2} \log(1/\delta) \right) = \tilde{O} \left( \frac{d_{\mathrm{LV}}^7 K^4 H^{12}}{\eta_{\min}^4} \log(1/\delta) \right),
$$

and the calls to LINEAR-FQI require $nHd_{\mathrm{LV}}$ samples in total. Therefore, the total sample complexity is

$$\tilde{O}\left(\max\left\{\frac{d_{\mathrm{LV}}^2 K^2 H^6}{\eta_{\min}^2}, \frac{1}{\varepsilon^2}\right\} \cdot Hd_{\mathrm{LV}}K^2 \log(|\Phi||\Upsilon|/\delta) + \frac{d_{\mathrm{LV}}^8 K^4 H^{13}}{\eta_{\min}^4}\log(1/\delta)\right).$$

# E   Maximum Likelihood Estimation

In this section we adapt classical results for maximum likelihood estimation in general parametric models. We consider a sequential conditional probability estimation setting with an instance space $\mathcal{X}$ and target space $\mathcal{Y}$ and with a conditional density $p(y \mid x) = f^\star(x, y)$. We are given a function class $\mathcal{F} : (\mathcal{X} \times \mathcal{Y}) \to \mathbb{R}$ with which to model the condition distribution $f^\star$, and we assume that $f^\star \in \mathcal{F}$, so that the problem is well-specified or realizable. We are given a dataset $D := \{(x_i, y_i)\}_{i=1}^n$, where $x_i \sim \mathcal{D}_i = \mathcal{D}_i(x_{1:i-1}, y_{1:i-1})$ and $y_i \sim p(\cdot \mid x_i)$. Note that $\mathcal{D}_i$ depends on the previous examples, so this is a martingale process. We optimize the maximum likelihood objective

$$\hat{f} := \operatorname*{argmax}_{f \in \mathcal{F}} \sum_{i=1}^n \log f(x_i, y_i). \tag{10}$$

The iid version of the following result is classical (c.f., Van de Geer, 2000, Chapter 7), but under-utilized in machine learning and reinforcement learning in particular. Our adaptation is inspired by Zhang (2006).

**Theorem 21.** *Fix $\delta \in (0, 1)$, assume $|\mathcal{F}| < \infty$ and $f^\star \in \mathcal{F}$. Then with probability at least $1 - \delta$*

$$\sum_{i=1}^n \mathbb{E}_{x \sim \mathcal{D}_i} \left\| \hat{f}(x, \cdot) - f^\star(x, \cdot) \right\|_{\mathrm{TV}}^2 \le 2\log(|\mathcal{F}|/\delta).$$

**Remark 22.** *Given a class of discriminators $\mathcal{G} : (\mathcal{X}, \mathcal{Y}) \mapsto [-1, 1]$, an alternative is to consider the following (conditional) "generative adversarial" objective:*

$$\hat{f} = \operatorname*{argmin}_{f \in \mathcal{F}} \max_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \left(g(x_i, y_i) - \mathbb{E}[g(x_i, y) \mid y \sim f(x, \cdot)]\right).$$

*This is the natural objective associated with the distance function induced by $\mathcal{G}$ (Arora et al., 2017), and is also related to other GAN-style approaches. Owing to the realizability assumption, $f^\star$ will always have low objective value, scaling with the complexity of $\mathcal{G}$. Additionally, if $\mathcal{G}$ is expressive enough, one can establish a guarantee similar to Theorem 21, which can then be used in the analysis of FLAMBE. Formally, a sufficient condition is that $\mathcal{G}$ contains the indicators of the Scheffe sets for all pairs $f, f' \in \mathcal{F}$, in which case the total variation guarantee can be obtained by standard uniform convergence arguments. See Devroye and Lugosi (2012); Sun et al. (2019) for more details.*

**Remark 23.** *We also remark that the proof of Theorem 21 actually establishes convergence in the squared Hellinger distance. We obtain the total variation guarantee simply by observing that the squared Hellinger distance dominates the squared total variation distance.*

We prove Theorem 21 in this section. We begin with a decoupling inequality. Let $D$ denote the dataset and let $D'$ denote a *tangent sequence* $\{(x_i', y_i')\}_{i=1}^n$ where $x_i' \sim \mathcal{D}_i(x_{1:i-1}, y_{1:i-1})$ and $y_i' \sim p(\cdot \mid x_i')$. Note here that $x_i'$ depends on the original sequence, and so the tangent sequence is independent conditional on $D$.

**Lemma 24.** *Let $D$ be a dataset of $n$ examples, and let $D'$ be a tangent sequence. Let $L(f, D) = \sum_{i=1}^n \ell(f, (x_i, y_i))$ be any function that decomposes additively across examples where $\ell$ is any function, and let $\hat{f}(D)$ be any estimator taking as input random variable $D$ and with range $\mathcal{F}$. Then*

$$\mathbb{E}_D\left[\exp\left(L(\hat{f}(D), D) - \log \mathbb{E}_{D'} \exp(L(\hat{f}(D), D')) - \log|\mathcal{F}|\right)\right] \le 1$$

38

Observe that in the second term, the "loss function" takes as input $D'$, but the estimator takes as input $D$. As such, the above inequality decouples the estimator from the loss.

*Proof.* Let $\pi$ be the uniform distribution over $\mathcal{F}$ and let $g : \mathcal{F} \to \mathbb{R}$ be any function. Define $\mu(f) := \frac{\exp(g(f))}{\sum_f \exp(g(f))}$, which is clearly a probability distribution. Now consider any other probability distribution $\hat{\pi}$ over $\mathcal{F}$:

$$0 \le \mathrm{KL}(\hat{\pi}||\mu) = \sum_f \hat{\pi}(f) \log(\hat{\pi}(f)) + \sum_f \hat{\pi}(f) \log\left(\sum_{f'} \exp(g(f'))\right) - \sum_f \hat{\pi}(f)g(f)$$

$$= \mathrm{KL}(\hat{\pi}||\pi) - \sum_f \hat{\pi}(f)g(f) + \log \mathbb{E}_{f \sim \pi} \exp(g(f))$$

$$\le \log|\mathcal{F}| - \sum_f \hat{\pi}(f)g(f) + \log \mathbb{E}_{f \sim \pi} \exp(g(f)).$$

Re-arranging, it holds that

$$\sum_f \hat{\pi}(f)g(f) - \log|\mathcal{F}| \le \log \mathbb{E}_{f \sim \pi} \exp(g(f)).$$

We instantiate this bound with $\hat{\pi} = \mathbf{1}\{\hat{f}(D)\}$ and $g(f) = L(f, D) - \log \mathbb{E}_{D'} \exp(L(f, D'))$ to obtain, for any $D$

$$L(\hat{f}(D), D) - \log \mathbb{E}_{D'} \exp(L(\hat{f}(D), D')) - \log|\mathcal{F}| \le \log \mathbb{E}_{f \sim \pi} \frac{\exp\left(L(f, D)\right)}{\mathbb{E}_{D'} \exp(L(f, D'))}.$$

Exponentiating both sides and then taking expectation over $D$, we obtain

$$\mathbb{E}_D \left[ \exp(L(\hat{f}(D), D) - \log \mathbb{E}_{D'} \exp(L(\hat{f}(D), D')) - \log|\mathcal{F}|) \right]$$

$$\le \mathbb{E}_{f \sim \pi} \mathbb{E}_D \frac{\exp\left(L(f, D)\right)}{\mathbb{E}_{D'}[\exp(L(f, D')) \mid D]} = 1.$$

The last equality follows since, conditional on $D$, the tangent sequence $D'$ is independent. Therefore,

$$\mathbb{E}_{D'} \left[ \exp(L(f, D')) \mid D \right] = \prod_{i=1}^n \mathbb{E}_{(x_i', y_i') \sim \mathcal{D}_i} \left[ \exp(\ell(f, (x_i', y_i'))) \right],$$

which allows us to peel off terms starting from $n$ down to $1$ and cancel them with those in the numerator. $\qquad\square$

The next lemma translates from TV-distance to a loss function that is closely related to the KL divergence.

**Lemma 25.** *For any two conditional probability densities $f_1, f_2$ and any distribution $\mathcal{D} \in \Delta(\mathcal{X})$ we have*

$$\mathbb{E}_{x \sim \mathcal{D}} \|f_1(x, \cdot) - f_2(x, \cdot)\|_{\mathrm{TV}}^2 \le -2 \log \mathbb{E}_{x \sim \mathcal{D}, y \sim f_2(\cdot|x)} \exp\left(-\frac{1}{2} \log(f_2(x, y)/f_1(x, y))\right)$$

*Proof.* Let us begin by relating the total variation distance, which appears on the left hand side, to the (squared) Hellinger distance, which for densities $p, q$ over a domain $\mathcal{Z}$ is defined as

$$\mathrm{H}^2(q||p) := \int \left(\sqrt{p(z)} - \sqrt{q(z)}\right)^2 dz.$$

Lemma 2.3 in Tsybakov (2008) asserts that

$$\|p(\cdot) - q(\cdot)\|_{\mathrm{TV}}^2 \leq \mathrm{H}^2(q||p) \cdot \left(1 - \frac{\mathrm{H}^2(q||p)}{4}\right) \leq \mathrm{H}^2(q||p),$$

where the final inequality uses non-negativity of the Hellinger distance. Next, note that we can also write

$$\mathrm{H}^2(q||p) = \int p(z) + q(z) - 2\sqrt{p(z)q(z)}dz = 2 \cdot \mathbb{E}_{z \sim q}\left[1 - \sqrt{p(z)/q(z)}\right]$$

$$\leq -2\log \mathbb{E}_{z \sim q}\sqrt{p(z)/q(z)} = -2\log \mathbb{E}_{z \sim q}\exp\left(-\frac{1}{2}\log(q(z)/p(z))\right).$$

Here the inequality follows from the fact that $1 - x \leq -\log(x)$. The result follows by applying this argument to $\mathbb{E}_{x \sim \mathcal{D}}\|f_1(x, \cdot) - f_2(x, \cdot)\|_{\mathrm{TV}}^2$. $\qquad\square$

*Proof of Theorem 21.* First note that Lemma 24 can be combined with the Chernoff method to obtain an exponential tail bound: with probability $1 - \delta$ we have

$$-\log \mathbb{E}_{D'}\exp(L(\hat{f}(D), D')) \leq -L(\hat{f}(D), D) + \log|\mathcal{F}| + \log(1/\delta).$$

Now we set $L(f, D) = \sum_{i=1}^{n} -1/2 \cdot \log(f^\star(x_i, y_i)/f(x_i, y_i))$ where $D$ is a dataset $\{(x_i, y_i)\}_{i=1}^{n}$ (and $D' = \{(x_i', y_i')\}_{i=1}^{n}$ is a tangent sequence). With this choice, the right hand side is

$$\sum_{i=1}^{n} \frac{1}{2}\log(f^\star(x_i, y_i)/\hat{f}(x_i, y_i)) + \log|\mathcal{F}| + \log(1/\delta) \leq \log|\mathcal{F}| + \log(1/\delta),$$

since $\hat{f}$ is the empirical maximum likelihood estimator and we are in the well-specified setting. On the other hand, the left hand side is

$$-\log \mathbb{E}_{D'}\left[\exp\left(\sum_{i=1}^{n} -1/2\log\left(\frac{f^\star(x_i', y_i')}{\hat{f}(x_i', y_i')}\right)\right) \mid D\right] = -\sum_{i=1}^{n}\log \mathbb{E}_{x,y \sim \mathcal{D}_i}\exp\left(-1/2\log\left(\frac{f^\star(x, y)}{\hat{f}(x, y)}\right)\right)$$

$$\geq \frac{1}{2}\sum_{i=1}^{n}\mathbb{E}_{x \sim \mathcal{D}_i}\left\|\hat{f}(x, \cdot) - f^\star(x, \cdot)\right\|_{\mathrm{TV}}^2.$$

Here the first identity uses the independence of the terms, which holds because $\hat{f}$ is independent of the dataset $D'$. The second inequality is Lemma 25. This yields the theorem. $\qquad\square$

# F    Auxilliary Lemmas

**Lemma 26** (Elliptical Potential Lemma). *Consider a sequence of $d \times d$ positive semidefinite matrices $X_1, \ldots, X_T$ with $\max_t \mathrm{tr}(X_t) \leq 1$ and define $M_0 = \lambda I_{d \times d}, \ldots, M_t = M_{t-1} + X_t$. Then*

$$\sum_{t=1}^{T}\mathrm{tr}(X_t M_{t-1}^{-1}) \leq (1 + 1/\lambda)d\log(1 + T/d).$$

*Proof.* Observe that by concavity of the $\log \det(\cdot)$ function, we have

$$\log(\det(M_{t-1})) \leq \log(\det(M_t)) + \mathrm{tr}(M_t^{-1}(M_{t-1} - M_t)).$$

40

Re-arranging and summing across all rounds $t$ yields

$$\sum_{t=1}^{T} \text{tr}(X_t M_t^{-1}) \leq \sum_{t=1}^{T} \log(\det(M_t)) - \log(\det(M_{t-1})) = \log(\det(M_T)) - d\lambda.$$

We will drop the negative term. By the spectral version of the AM-GM inequality and linearity of trace, we upper bound the last term:

$$\det(M_T)^{1/d} \leq \text{tr}(M_T)/d \leq 1 + T/d.$$

Now, we must convert from $M_t^{-1}$ to $M_{t-1}^{-1}$ on the left hand side. Fix a round $t$ and let us write $X_t = VV^\top$, which is always possible as $X_t$ is positive semidefinite. Then by the Woodbury identity

$$\begin{aligned}
\text{tr}(X_t M_t^{-1}) &= \text{tr}\left(V^\top (M_{t-1} + VV^\top)^{-1} V\right) \\
&= \text{tr}(V^\top M_{t-1}^{-1} V) - \text{tr}(V^\top M_{t-1}^{-1} V (I + V^\top M_{t-1}^{-1} V)^{-1} V^\top M_{t-1}^{-1} V).
\end{aligned}$$

All matrices are simultaneously diagonalizable, so we may pass to a common eigendecomposition. In particular, with the eigendecomposition $V^\top M_{t-1}^{-1} V = \sum_{i=1}^{d} \lambda_i u_i u_i^\top$, we obtain

$$\text{tr}(X_t M_t^{-1}) = \sum_{i=1}^{d} \lambda_i - \frac{\lambda_i^2}{1 + \lambda_i} = \sum_{i=1}^{d} \frac{\lambda_i}{1 + \lambda_i} \geq \frac{1}{1 + 1/\lambda} \sum_{i=1}^{d} \lambda_i = \frac{1}{1 + 1/\lambda} \text{tr}(X_t M_{t-1}^{-1}).$$

The inequality follows from the fact that $\lambda_i \leq \left\|V^\top M_{t-1}^{-1} V\right\|_2 \leq 1/\lambda$ due to our initial conditions on $M_0$ and the normalization for $X_t$. $\qquad\square$

**Corollary 27.** *Consider the setup of [Lemma 26](#) and further assume that for each $t$, we have $\text{tr}(X_t M_{t-1}^{-1}) \geq \beta > 0$. Then $T \leq 2(1 + 1/\lambda)d \log(1 + 2(1 + 1/\lambda)/\beta)/\beta$.*

*Proof.* The stated assumption and [Lemma 26](#) implies that $T \leq (1 + 1/\lambda)d \log(1 + T/d)/\beta$. We claim that if $T \leq 2(1 + 1/\lambda)d \log(1 + 2(1 + 1/\lambda)/\beta)/\beta$ then a weakening of this bound is

$$\begin{aligned}
T &\leq \frac{(1 + 1/\lambda)d}{\beta} \log(1 + T/d)/\beta \leq \frac{(1 + 1/\lambda)d}{\beta} \log\left(1 + \frac{2(1 + 1/\lambda)\log(1 + 2(1 + 1/\lambda)/\beta)}{\beta}\right) \\
&\leq \frac{(1 + 1/\lambda)d}{\beta} \log\left(1 + \left(\frac{2(1 + 1/\lambda)}{\beta}\right)^2\right) \leq \frac{2(1 + 1/\lambda)d}{\beta} \log\left(1 + \frac{2(1 + 1/\lambda)}{\beta}\right).
\end{aligned}$$

Therefore, we have established an upper bound on $T$. $\qquad\square$

# References

Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, 2014.

Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *International Conference on Machine Learning*, 2017.

Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. In *International Conference on Machine Learning*, 2019.

Andre Barreto, Doina Precup, and Joelle Pineau. Reinforcement learning using kernel-based stochastic factorization. In *Advances in Neural Information Processing Systems*, 2011.

André MS Barreto and Marcelo D Fragoso. Computing the stationary distribution of a finite markov chain through stochastic factorization. *SIAM Journal on Matrix Analysis and Applications*, 2011.

Jonathan Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 2000.

Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013.

Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. *arXiv:1912.05830*, 2019.

Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, 2019.

Joel E Cohen and Uriel G Rothblum. Nonnegative ranks, decompositions, and factorizations of nonnegative matices. *Linear Algebra and its Applications*, 1993.

Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *Conference on Learning Theory*, 2008.

Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2012.

Kefan Dong, Jian Peng, Yining Wang, and Yuan Zhou. $\sqrt{n}$-regret for learning in Markov decision processes with function approximation and low Bellman rank. *arXiv:1909.02506*, 2019.

Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *International Conference on Learning Representations*, 2019a.

Simon S Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudík, and John Langford. Provably efficient RL with rich observations via latent state decoding. In *International Conference on Machine Learning*, 2019b.

Simon S Du, Yuping Luo, Ruosong Wang, and Hanrui Zhang. Provably efficient Q-learning with function approximation via distribution shift error checking oracle. In *Advances in Neural Information Processing Systems*, 2019c.

Simon S Du, Jason D Lee, Gaurav Mahajan, and Ruosong Wang. Agnostic Q-learning with function approximation in deterministic systems: Tight bounds on approximation error and sample complexity. *arXiv:2002.07125*, 2020.

Yaqi Duan, Mengdi Wang, Zaiwen Wen, and Yaxiang Yuan. Adaptive low-nonnegative-rank approximation for state aggregation of Markov chains. *SIAM Journal on Matrix Analysis and Applications*, 2020.

Jack Edmonds. Maximum matching and a polyhedron with 0,1-vertices. *Journal of Research of the National Bureau of Standards–B*, 1965.

Fei Feng, Ruosong Wang, Wotao Yin, Simon S Du, and Lin F Yang. Provably efficient exploration for rl with unsupervised learning. *arXiv:2003.06898*, 2020.

Mikhail Figurnov, Shakir Mohamed, and Andriy Mnih. Implicit reparameterization gradients. In *Advances in Neural Information Processing Systems*, 2018.

Samuel Fiorini, Volker Kaibel, Kanstantsin Pashkovich, and Dirk Oliver Theis. Combinatorial bounds on nonnegative rank and extended formulations. *Discrete Mathematics*, 2013.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, 2014.

Elad Hazan, Sham M Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy exploration. In *International Conference on Machine Learning*, 2019.

Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparametrization with gumbel-softmax. In *International Conference on Learning Representations*, 2017.

Nan Jiang, Alex Kulesza, and Satinder Singh. Abstraction selection in model-based reinforcement learning. In *International Conference on Machine Learning*, 2015.

Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low Bellman rank are PAC-learnable. In *International Conference on Machine Learning*. JMLR. org, 2017.

Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv:1907.05388*, 2019.

Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. *arXiv:2002.02794*, 2020.

Tor Lattimore and Csaba Szepesvari. Learning with good feature representations in bandits and in rl with a generative model. *arXiv:1911.07676*, 2019.

Tor Lattimore, Marcus Hutter, Peter Sunehag, et al. The sample-complexity of general reinforcement learning. In *International Conference on Machine Learning*. Journal of Machine Learning Research, 2013.

Michael L Littman and Richard S Sutton. Predictive representations of state. In *Advances in Neural Information Processing Systems*, 2002.

Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 2016.

Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. *arXiv:1911.05815*, 2019.

Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *Conference on Artificial Intelligence and Statistics*, 2020.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv:1807.03748*, 2018.

Ronald Ortner, Odalric-Ambrym Maillard, and Daniil Ryabko. Selecting near-optimal approximate state representations in reinforcement learning. In *International Conference on Algorithmic Learning Theory*. Springer, 2014.

Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In *Advances in Neural Information Processing Systems*, 2014.

Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017.

Steffen Rendle, Christoph Freudenthaler, and Lars Schmidt-Thieme. Factorizing personalized Markov chains for next-basket recommendation. In *International Conference on World Wide Web*, 2010.

Thomas Rothvoß. The matching polytope has exponential extension complexity. *Journal of the ACM*, 2017.

Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, 2013.

Satinder Singh, Michael R James, and Matthew R Rudary. Predictive state representations: a new theory for modeling dynamical systems. In *Conference on Uncertainty in Artificial Intelligence*, 2004.

Aravind Srinivas, Michael Laskin, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv:2004.04136*, 2020.

Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based RL in contextual decision processes: PAC bounds and exponential improvements over model-free approaches. In *Conference on Learning Theory*, 2019.

Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schulman, Filip DeTurck, and Pieter Abbeel. #Exploration: A study of count-based exploration for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2017.

Michael Thon and Herbert Jaeger. Links between multiplicity automata, observable operator models and predictive state representations: a unified learning framework. *The Journal of Machine Learning Research*, 2015.

Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive estimation reveals topic posterior information to linear models. *arXiv:2003.02234*, 2020.

Alexandre B Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.

Sara Van de Geer. *Empirical Processes in M-estimation*. Cambridge University Press, 2000.

Benjamin Van Roy and Shi Dong. Comments on the Du-Kakade-Wang-Yang lower bounds. *arXiv:1911.07910*, 2019.

Yining Wang, Ruosong Wang, Simon S Du, and Akshay Krishnamurthy. Optimism in reinforcement learning with generalized linear function approximation. *arXiv:1912.04136*, 2019.

Max Welling and Yee W Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *International Conference on Machine Learning*, 2011.

Zheng Wen and Benjamin Van Roy. Efficient exploration and value function generalization in deterministic systems. In *Advances in Neural Information Processing Systems*, 2013.

Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive features. In *International Conference on Machine Learning*, 2019a.

Lin F Yang and Mengdi Wang. Reinforcement leaning in feature space: Matrix bandit, kernels, and regret bound. *arXiv:1905.10389*, 2019b.

Mihalis Yannakakis. Expressing combinatorial optimization problems by linear programs. *Journal of Computer and System Sciences*, 1991.

Hengshuai Yao, Csaba Szepesvári, Bernardo Avila Pires, and Xinhua Zhang. Pseudo-mdps and factored linear action models. In *IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*, 2014.

Tong Zhang. From $\epsilon$-entropy to KL-entropy: Analysis of minimum information complexity density estimation. *The Annals of Statistics*, 2006.