

---

# A maximum-entropy approach to off-policy evaluation in average-reward MDPs

---

Nevena Lazić\*   Dong Yin\*   Mehrdad Farajtabar\*   Nir Levine\*   Dilan Görür\*  
 Chris Harris†   Dale Schuurmans†

## Abstract

This work focuses on off-policy evaluation (OPE) with function approximation in infinite-horizon undiscounted Markov decision processes (MDPs). For MDPs that are ergodic and linear (i.e. where rewards and dynamics are linear in some known features), we provide the first finite-sample OPE error bound, extending existing results beyond the episodic and discounted cases. In a more general setting, when the feature dynamics are approximately linear and for arbitrary rewards, we propose a new approach for estimating stationary distributions with function approximation. We formulate this problem as finding the maximum-entropy distribution subject to matching feature expectations under empirical dynamics. We show that this results in an exponential-family distribution whose sufficient statistics are the features, paralleling maximum-entropy approaches in supervised learning. We demonstrate the effectiveness of the proposed OPE approaches in multiple environments.

## 1 Introduction

Recently, there have been considerable advances in reinforcement learning (RL), with algorithms achieving impressive performance on game playing and simple robotic tasks. Successful approaches typically learn through direct (online) interaction with the environment. However, in many real applications, access to the environment is limited to a fixed dataset, due to considerations of cost, safety, or time. One key challenge in this setting is off-policy evaluation (OPE): the task of evaluating the performance of a target policy given samples collected by a behavior policy.

The focus of our work is OPE in infinite-horizon undiscounted MDPs, which capture long-horizon tasks such as game playing, routing, and the control of physical systems. Most recent state-of-the-art OPE methods for this setting estimate the ratios of stationary distributions of the target and behavior policy [Liu et al., 2018, Nachum et al., 2019a, Wen et al., 2020, Nachum and Dai, 2020]. These approaches typically produce estimators that are consistent, but have no finite-sample guarantees, and even the existing guarantees may not hold with function approximation. One exception is the recent work of Duan and Wang [2020], which relies on linear function approximation. They assume that the MDP is linear (i.e. that rewards and dynamics are linear in some known feature space) and analyze OPE in episodic and discounted MDPs when given a fixed dataset of i.i.d. trajectories. They establish a finite-sample instance-dependent error upper bound for regression-based fitted Q-iteration (FQI), and a nearly-matching minimax lower bound.

Our work extends the results of Duan and Wang [2020] to the setting of undiscounted ergodic linear MDPs and non-i.i.d. data (coming from a single trajectory). We provide the first finite-sample OPE error bound for this case; our bound scales similarly to that of Duan and Wang [2020], but depends

---

\*DeepMind

†Google

on the MDP mixing time rather than horizon or discount. We are not aware of any similar results for off-policy evaluation in average-reward MDPs. Indeed, while OPE with linear function approximation has been well-studied for discounted MDPs [Geist and Scherrer, 2014, Dann et al., 2014, Yu, 2010a], in the undiscounted setting even showing convergence of standard methods presents some difficulties (see the discussion in Yu [2010b] for more details).

Beyond linear MDPs, we consider MDPs in which rewards are non-linear, while the state-action dynamics are still (approximately) linear in some features. Here we propose a novel approach for estimating stationary distributions with function approximation: we maximize the distribution entropy subject to matching feature expectations under the empirical dynamics. Interestingly, this results in an exponential family distribution whose sufficient statistics are the features, paralleling the well-known maximum entropy approach to supervised learning [Jaakkola et al., 2000]. We demonstrate the effectiveness of our proposed OPE approaches in multiple environments.

## 2 Preliminaries

**Problem definition.** We are interested in learning from batch data in infinite-horizon ergodic Markov decision processes (MDPs). An MDP is a tuple  $(\mathcal{S}, \mathcal{A}, r, P)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function, and  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  is the transition probability function. For ease of exposition, we will assume that states and actions are discrete, but similar ideas apply to continuous state and action spaces. A policy  $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  is a mapping from a state to a distribution over actions. We will use  $\mathbb{I}_{\pi}$  to denote the transition kernel from a state-action pair  $(s, a)$  to the next pair  $(s', a')$  under  $\pi$ . In an ergodic MDP, every policy induces a single recurrent class of states, i.e. any state can be reached from any other state. The expected average reward of a policy is defined as

$$J_{\pi} = \lim_{T \rightarrow \infty} \mathbf{E} \left[ \frac{1}{T} \sum_{t=1}^T r(s_t, a_t) \right] \quad \text{where } s_{t+1} \sim P(\cdot | s_t, a_t) \text{ and } a_t \sim \pi(\cdot | s_t).$$

Assume we are given a trajectory of  $T$  transitions  $\mathcal{D}_{\beta} = \{(s_t, a_t, r_t)\}_{t=1}^{T+1}$  generated by a behavior policy  $\beta$  in an unknown MDP. The *off-policy evaluation* problem is the task of estimating  $J_{\pi}$  for a target policy  $\pi$ .

**Stationary distributions.** Let  $\mu_{\pi}(s)$  be the stationary state distribution of a policy  $\pi$ , and let  $d_{\pi}(s, a) = \mu_{\pi}(s)\pi(a|s)$  be the stationary state-action distribution. These distributions satisfy the flow constraint

$$\mu_{\pi}(s') = \sum_{a'} d_{\pi}(s', a') = \sum_{s, a} d_{\pi}(s, a) P(s' | s, a). \quad (1)$$

The expected average reward can equivalently be written as  $J_{\pi} = \mathbf{E}_{(s, a) \sim d_{\pi}}[r(s, a)]$ . Thus, one approach to learning in MDPs from batch data involves estimating or optimizing  $d_{\pi}$  subject to (1). In particular, given data sampled from  $d_{\beta}$  and distribution estimates  $\hat{d}_{\pi}$  and  $\hat{d}_{\beta}$ , we can estimate  $J_{\pi}$  as  $\hat{J}_{\pi} = \frac{1}{T} \sum_{t=1}^T \frac{\hat{d}_{\pi}(s_t, a_t)}{\hat{d}_{\beta}(s_t, a_t)} r_t$ , as proposed by Liu et al. [2018].

**Linear MDPs.** When the state-action space is large or continuous-valued, a common approach to evaluating or optimizing a policy is to use function approximation. Define the conditional transition operator  $\mathcal{P}^{\pi}$  of a policy  $\pi$  as

$$\mathcal{P}^{\pi} f(s, a) := \mathbf{E}_{s' \sim P(\cdot | s, a), a' \sim \pi(\cdot | s')} [f(s', a') | s, a]. \quad (2)$$

With function approximation, it is convenient to assume that for any policy,  $\mathcal{P}^{\pi}$  operates within a particular function class  $\mathcal{F}$ , i.e. for any  $f \in \mathcal{F}$ ,  $\mathcal{P}^{\pi} f \in \mathcal{F}$  [Duan and Wang, 2020]. We will assume that  $\mathcal{F}$  is the set of functions linear in some (known or pre-learned) features  $\phi(s, a) \in \mathbb{R}^m$ , such that for some matrix  $M_{\pi} \in \mathbb{R}^{m \times m}$ ,

$$\mathcal{P}^{\pi} \phi(s, a) = \sum_{s', a'} \pi(a' | s') P(s' | s, a) \phi(s', a')^{\top} = \phi(s, a)^{\top} M_{\pi} + b_{\pi}^{\top}. \quad (3)$$

Note that, unlike existing work, we specifically include a bias term  $b_{\pi}$  in the above model. When  $\phi(s, a)$  is a binary indicator vector for  $(s, a)$ ,  $M_{\pi}$  corresponds to the state-action transition matrix and  $b_{\pi} = 0$ . However,  $b_{\pi}$  is non-zero in other settings, such as MDPs with linear-Gaussian dynamics. Similarly to Duan and Wang [2020], we will assume that rewards  $r(s, a)$  are linear in the same features:  $r(s, a) = \phi(s, a)^{\top} w$ . This assumption will be required for the purpose of analysis.

### 3 Off-policy evaluation

#### 3.1 Maximum-entropy stationary distribution estimation

Given a policy  $\pi(a|s)$ , in order to compute an off-policy estimate of  $J_\pi$ , we only need to estimate the stationary state distribution  $\mu_\pi(s)$ . We formulate this as a maximum-entropy problem subject to matching feature expectations:

$$\min_{\mu \in \Delta_S} \sum_s \mu(s) \ln \mu(s) \quad (4)$$

$$\text{s.t.} \quad \sum_{s', a'} \mu_\pi(s') \pi(a'|s') \phi(s', a') = \sum_{s, a} \mu_\pi(s) \pi(a|s) \sum_{s', a'} P(s'|s, a) \pi(a'|s') \phi(s', a'). \quad (5)$$

Note that we have relaxed the original flow constraint (1) over all state-action pairs to only require feature expectations to match, similarly to the maximum-entropy principle for supervised learning [Jaakkola et al., 2000]. Furthermore, under the linear MDP assumption and given the model parameters  $(M_\pi, b_\pi)$ , the feature expectation constraint can be written as

$$\sum_s \mu_\pi(s) \phi(s, \pi)^\top (I - M_\pi) = b_\pi^\top, \quad (6)$$

where  $\phi(s, \pi) = \sum_a \pi(a|s) \phi(s, a)$  are feature expectations under the policy. In Appendix A, we show that the optimal solution is an exponential-family distribution of the following form:

$$\mu_\pi(s|\theta_\pi, M_\pi) = \exp(\phi(s, \pi)^\top (I - M_\pi) \theta_\pi - F(\theta_\pi|M_\pi)) \quad (7)$$

where  $F(\theta_\pi|M_\pi)$  is the log-partition function. The parameters  $\theta_\pi$  are the solution of the dual problem:

$$\theta_\pi = \arg \min_{\theta} D(\theta) := F(\theta|M_\pi) - \theta^\top b_\pi. \quad (8)$$

Note that the dual is convex, due to the convexity of the log-partition function in exponential families. Given a batch of data, we estimate the stationary distribution  $\mu_\pi$  by first estimating  $\widehat{M}_\pi$  and  $\widehat{b}_\pi$  using linear regression (see (11)), and then computing a parameter estimate as  $\widehat{\theta}_\pi = \arg \min_{\theta} F(\theta|\widehat{M}_\pi) - \widehat{b}_\pi^\top \theta$ . When the log-partition function  $F(\theta|\widehat{M}_\pi)$  is intractable, we can optimize the dual using stochastic gradient descent. Noting that  $\nabla_{\theta} F(\theta|M_\pi) = \mathbf{E}_{\mu_\pi}[(I - M_\pi^\top) \phi(s, \pi)]$ , we can obtain an (almost) unbiased gradient estimate using importance weights:

$$\widehat{\nabla}_{\theta} F(\theta|\widehat{M}_\pi) \propto \sum_{s \in \mathcal{D}_\beta} \frac{\hat{\mu}_\pi(s|\theta, \widehat{M}_\pi)}{\hat{\mu}_\beta(s|\widehat{\theta}_\beta, \widehat{M}_\beta)} (I - \widehat{M}_\pi^\top) \phi(s, \pi) \quad (9)$$

where  $\hat{\mu}_\beta(s|\widehat{\theta}_\beta, \widehat{M}_\beta)$  is an estimate of the stationary distribution of the behavior policy computed using the same approach (we assume that the behavior policy is known and otherwise estimate it from the data). Finally, we evaluate the policy as

$$\widehat{J}_\pi = \sum_{t=1}^T \rho_t r_t \quad \text{where} \quad \rho_t = \frac{\hat{\mu}_\pi(s_t) \pi(a_t|s_t)}{\hat{\mu}_\beta(s_t) \beta(a_t|s_t)}.$$

In practice, it may be beneficial to normalize the distribution weights  $\rho_t$  to sum to 1, known as weighted importance sampling [Rubinstein, 1981, Koller and Friedman, 2009, Mahmood et al., 2014]. This results in an estimate that is biased but consistent, and often of much lower variance; the same technique can be applied to the gradient weights following Chen and Luss [2018]. When the log-normalizing constant is intractable, we can normalize the distributions empirically.

**Linear rewards.** When the rewards are linear in the features,  $r(s, a) = \phi(s, a)^\top w$ , and  $b_\pi \neq \mathbf{0}$ , there is a faster way to estimate  $J_\pi$ . Noting that since  $J_\pi = \sum_{s, a} d_\pi(s, a) \phi(s, a)^\top w$ , we only need to estimate  $e_\pi := \sum_{s, a} d_\pi(s, a) \phi(s, a)$  rather than the full distribution. Under the linear MDP assumption,  $e_\pi^\top = b_\pi^\top (I - M_\pi)^{-1}$ . Thus, given estimates of the model and reward parameters  $\widehat{M}_\pi, \widehat{b}_\pi, \widehat{w}$ , we can evaluate the policy as

$$\widehat{J}_\pi = \widehat{b}_\pi^\top (I - \widehat{M}_\pi)^{-1} \widehat{w}. \quad (10)$$

### 3.2 OPE error analysis.

Our analysis requires the following assumptions.

**Assumption A1 (Mixing coefficient)** There exists a constant  $\kappa > 0$  such that for any state-action distribution  $d$ ,

$$\|(d_\beta - d)^\top \Pi_\beta\|_1 \leq \exp(-1/\kappa) \|d_\beta - d\|_1$$

where  $\Pi_\beta$  is the transition matrix from  $(s, a)$  to  $(s', a')$  under the policy  $\beta$ .

**Assumption A2 (Bounded linearly independent features)** Let  $\bar{\phi}(s, a)^\top := [\phi(s, a)^\top \ 1]$ . We assume that  $\max_{s,a} \|\bar{\phi}(s, a)\|_2 \leq C_\Phi$  for some constant  $C_\Phi$ . Let  $\Phi$  be an  $|\mathcal{S}||\mathcal{A}| \times (m+1)$  matrix whose rows are feature vectors  $\bar{\phi}(s, a)$ . We assume that the columns of  $\Phi$  are linearly independent.

**Assumption A3 (Feature excitation)** For a policy  $\pi$  with stationary distribution  $d_\pi(s, a)$ , define  $\Sigma_\pi = \mathbf{E}_{(s,a) \sim d_\pi} [\bar{\phi}(s, a) \bar{\phi}(s, a)^\top]$ . We assume that  $\lambda_{\min}(\Sigma_\beta) \geq \sigma > 0$  and  $\lambda_{\min}(\Sigma_\pi) \geq \sigma_\pi > 0$ .

The above assumptions mean that the exploration policy  $\beta(a|s)$  mixes fast and is exploratory, in the sense that the stationary distribution spans all dimensions of the feature space. These assumptions allow us to bound the model error. We also require the evaluated policy to span the feature space for somewhat technical reasons, in order to bound the policy evaluation error.

Assume that rewards are linear in the features,  $r(s, a) = \phi(s, a)^\top w$ . Given a trajectory  $\{(s_t, a_t, r_t)\}_{t=1}^{T+1}$ , we estimate  $M_\pi$ ,  $b_\pi$ , and  $w$  using regularized least squares:

$$\begin{bmatrix} \widehat{M}_\pi \\ \widehat{b}_\pi^\top \end{bmatrix} = \left( \Lambda + \sum_{t=1}^T \bar{\phi}(s_t, a_t) \bar{\phi}(s_t, a_t)^\top \right)^{-1} \sum_{t=1}^T \bar{\phi}(s_t, a_t) \phi(s_{t+1}, \pi)^\top \quad (11)$$

$$\widehat{w} = \left( \sum_{t=1}^T \phi(s_t, a_t) \phi(s_t, a_t)^\top \right)^{-1} \sum_{t=1}^T \phi(s_t, a_t) r_t \quad (12)$$

where  $\Lambda$  is a regularizer and  $\bar{\phi}(s, a) = [\phi(s, a)^\top \ 1]$ . For the purpose of simplifying the analysis, we let  $\Lambda = \alpha \sum_{t=1}^T \bar{\phi}(s_t, a_t) \bar{\phi}(s_t, a_t)^\top$ ; in practice it may be better to use a diagonal matrix. Let  $W_\pi = \begin{bmatrix} M_\pi & \mathbf{0} \\ b_\pi^\top & 1 \end{bmatrix}$  and similarly  $\widehat{W}_\pi = \begin{bmatrix} \widehat{M}_\pi & \mathbf{0} \\ \widehat{b}_\pi^\top & 1 \end{bmatrix}$ . The following Lemma (proven in Appendix B) bounds the estimation error under Assumptions A1 and A3 for single-trajectory data:

**Lemma 3.1.** *Let assumptions A1, A2, and A3 hold, and let  $\alpha = C_\Phi^2 \sigma^{-1} \kappa / \sqrt{T}$ . Then with probability at least  $1 - \delta$ , for constants  $C$  and  $C_w$ ,*

$$\begin{aligned} \|\widehat{W}_\pi - W_\pi\|_2 &\leq C C_\Phi^4 \kappa \sigma^{-2} \sqrt{2 \ln(2(m+1)/\delta)/T} \\ \|w - \widehat{w}\|_2 &\leq C_w C_\Phi^2 \kappa \sigma^{-2} \sqrt{2 \ln(2m/\delta)/T} \|w\|_2. \end{aligned}$$

The following theorem bounds the policy evaluation error.

**Theorem 3.2 (Policy evaluation error).** *Let assumptions A1, A2, and A3 hold and assume that problem (4)-(5) is feasible. Then, for a constant  $C_J$ , with probability at least  $1 - \delta$ , the batch policy evaluation error is bounded as*

$$|J_\pi - \widehat{J}_\pi| \leq C_J C_\Phi^4 \kappa \sigma_\pi^{-1/2} \sigma^{-2} (1 + \alpha)^2 \sqrt{2 \ln(2(m+1)/\delta)/T} \|w\|_2. \quad (13)$$

The proof is given in Appendix C and relies on expressing evaluation error in terms of the model error, as well as on the contraction properties of the matrix  $(1 + \alpha)^{-1} W_\pi$ . While we do not provide a lower bound, note that the error scales similarly to the results of Duan and Wang [2020] for discounted MDPs, which nearly match the corresponding lower bound.

**Remark 1.** Theorem 3.2 holds for any feasible solution  $\mu$  of (4) and not necessarily just for the maximum-entropy distribution.

**Remark 2.** Our results are shown for the case of discrete states and actions and bounded-norm features. In the continuous case, similar conclusions would follow by arguments on the concentration and boundedness of  $\Sigma_\beta$  and  $\mathbf{E}_{d_{\beta,P}} [\bar{\phi}(s, a) \bar{\phi}(s', \pi)^\top]$ .

### 3.3 Policy improvement

The previous sections provides an approach for estimating the average reward, but it is unclear how to perform policy optimization. One possible formulation is to maximize the entropy-regularized expected reward:

$$\max_{d_\pi \in \Delta^{S \times A}} \sum_{s,a} d_\pi(s,a) (r(s,a) - \tau \ln d_\pi(s,a)) \quad \text{s.t.} \quad \sum_{s,a} d(s,a) \phi(s,a)^\top (I - \widehat{M}_\pi) = \widehat{b}_\pi.$$

Unfortunately, this is no longer a convex problem, as the model  $(\widehat{M}_\pi, \widehat{b}_\pi)$  depends on the optimization variables through  $\pi(a|s) = \frac{d(s,a)}{\sum_{a'} d(s,a')}$ . We describe some ways around this in Appendix F.

An alternative formulation is to construct a critic for the purpose of policy improvement. Let  $Q_\pi(s, a)$  be the state-action value function of a policy, corresponding to the value of taking action  $a$  in state  $s$  and then following the policy forever, and satisfying the Bellman equation

$$Q_\pi(s, a) + J_\pi = r(s, a) + \sum_{s', a'} P(s'|s, a) \pi(a'|s') Q_\pi(s', a') \quad (14)$$

The true  $Q_\pi$  and  $J_\pi$  minimize the Bellman error:

$$L_{BE}(Q, J) = \mathbf{E}_{(s,a) \sim d_\pi} [(Q(s, a) + J - r(s, a) - \mathbf{E}_{s' \sim P(\cdot|s,a), a' \sim \pi(\cdot|s')} [Q(s', a')])^2]. \quad (15)$$

In the off-policy case, the above expectation can be taken with respect to the stationary distribution of the behavior policy instead and importance-corrected if possible. Typically, the Bellman error cannot be minimized directly, as it includes an expectation over unknown  $P$ , and we only have access to a sample trajectory. However, under the linear MDP assumption, all  $Q$ -functions are linear, and thus we can estimate these expectations using the feature-dynamics model. Thus we have the following objective for fitting a critic with linear function approximation  $Q_\pi(s, a) = \phi(s, a)^\top v_\pi$ :

$$L_{MBE}(v, J | \widehat{M}_\pi, \widehat{b}_\pi) = \mathbf{E}_{(s,a) \sim d_\pi} [(\phi(s, a)^\top (I - \widehat{M}_\pi) v - \widehat{b}_\pi^\top v + J - r(s, a))^2] \quad (16)$$

Consider instead minimizing the absolute average Bellman error  $|L_{AVG}(v, J)|$ , where

$$L_{AVG}(v, J) := \mathbf{E}_{(s,a) \sim d_\pi} [\phi(s, a)^\top (I - \widehat{M}_\pi) v - \widehat{b}_\pi^\top v + J - r(s, a)]. \quad (17)$$

This error is minimized when  $J = \mathbf{E}_{d_\pi} [r(s, a)]$  and  $\mathbf{E}_{d_\pi} [\phi(s, a)^\top (I - \widehat{M}_\pi) v] = \widehat{b}_\pi^\top v$ . The second condition corresponds to driving the gradient of our dual objective to zero. The recent work of Xie and Jiang [2020] suggests that there are some theoretical advantages to minimizing the average rather than squared Bellman error in discounted MDPs. However, in the undiscounted case, it is unclear whether  $v$  is useful for policy improvement, as the minimizer is not a function of the reward. We leave further investigation of policy improvement in average-reward MDPs for future work.

## 4 Related work

The linear MDP assumption along with linear rewards implies that all value functions are linear. Thus we first discuss similarities between our approach and common linear action-value function methods in literature, and then give an broader overview of other related work.

**TD error.** The residual gradient algorithm of Baird [1995] minimizes the mean squared temporal difference (TD) error:

$$L_{TD}(v, J) = \frac{1}{T} \sum_{t=1}^T (\phi(s_t, a_t)^\top v + J - r_t - \phi(s'_t, \pi)^\top v)^2 \quad (18)$$

It is well-known that this objective is a biased and inconsistent estimate of the true Bellman error [Bradtke and Barto, 1996]. Correcting the bias requires double samples (two independent samples of  $s'$  for the same  $(s, a)$  pair), which may not be available in a single-trajectory dataset. More recent methods rely on fixed point iterations, using the parameters from the previous iteration to construct regression targets. In our setting, fitted Q-iteration (FQI) can be written as

$$v^{(k+1)}, J^{(k+1)} = \arg \min_{v, J} \sum_{t=1}^T (\phi(s_t, a_t)^\top v + J - r_t - \phi(s'_t, \pi)^\top v^{(k)})^2 \quad (19)$$

The convergence of FQI is guaranteed only in restricted cases (e.g. Antos et al. [2008]), and no guarantees exist in the undiscounted setting to the best of authors’ knowledge.

**PBE error.** Another class of methods minimize the projected Bellman error, which corresponds to only the error representable by the features. The advantage of this approach is that the error due to not having exact dynamics expectations is not correlated with the TD error [Sutton et al., 2009]. Let  $D_\beta = \text{diag}(d_\beta)$ , and let  $Q_\pi$  be the action-value function of  $\pi$  as a vector. In matrix form, the projected Bellman equation (PBE) can be written as

$$Q_\pi = G_\beta(r - J_\pi \mathbf{1} + \Pi_\pi Q_\pi), \quad \text{where } G_\beta = \Phi(\Phi^\top D_\beta \Phi)^{-1} \Phi^\top D_\beta \quad (20)$$

where  $\Phi$  excludes bias. Methods that attempt to solve (the sample-based version of) the above equation using fixed-point iteration may diverge if the matrix  $G_\beta \Pi_\pi$  is not contractive (has spectral radius greater than or equal to 1). The contractiveness condition has been shown to hold in the on-policy setting ( $\beta = \pi$ ) under mild assumptions by Yu and Bertsekas [2009], but need not hold in general. Practical solutions may ensure convergence by regularizing  $G_\beta$  as  $G_\beta^b = \Phi(\Phi^\top D_\beta \Phi + bI)^{-1} \Phi^\top D_\beta$ ; however, the required bias  $b$  may be large. Another approach, popular in the discounted setting, is to minimize the PBE error using least squares methods like LSTD and LSPE [Yu, 2010a, Dann et al., 2014, Geist and Scherrer, 2014]. A general complication in the average-reward case is that, unlike with discounted methods which only estimate  $Q_\pi$ , we also need to solve for  $J_\pi$ . One possible heuristic is to initially use a guess for  $J_\pi$ . However, the resulting projected equation may not have a solution [Yu, 2010a, Puterman, 2014], and for OPE,  $J_\pi$  is actually the quantity we are after.

**DualDICE** [Nachum and Dai, 2020]. DualDICE and related methods optimize the Lagrangian of the following linear programming (LP) formulation of the Bellman equation:

$$\min_{J, Q} -J \quad \text{s.t. } Q + J\mathbf{1} \leq r + \Pi_\pi Q.$$

For the average reward case, linear function approximation  $Q = \Phi v$ , and a linear MDP satisfying  $\Pi_\pi \Phi = \Phi W_\pi$ , the problem Lagrangian is (see also Zhang et al. [2020], Nachum and Dai [2020]):

$$\max_d \min_{J, v} L(d, J, v) = -J + d^\top (\Phi(I - W_\pi)v + J). \quad (21)$$

We solve an entropy-regularized version of the dual LP using a learned model  $\widehat{W}_\pi$ . This can be seen as a particular convex instantiation of DualDICE with function approximation - a linear feature model, linear value functions, and exponential-family stationary distributions.

**Dual approaches to batch RL.** Many recent methods for batch policy evaluation and optimization rely on estimating stationary distribution ratios that (approximately) respect the MDP dynamics [Liu et al., 2018, Nachum et al., 2019a,b, Wen et al., 2020]. In particular, Liu et al. [2018] impose a similar constraint to ours on matching feature expectations. However, while we enforce the constraint for a particular feature representation, they minimize the squared error of violating the constraint while maximizing over smooth feature functions in a reproducing kernel Hilbert space. Feng et al. [2019] minimize a kernel loss for solving the Bellman equation. We note that our approach can also be kernelized, by using kernel ridge regression in place of linear regression for the model. Most of the existing approaches yield consistent estimators, but have no finite-sample guarantees. One exception is the work of Duan and Wang [2020], which provides a minimax lower bound and nearly-matching finite-sample error bound in linear finite-horizon and discounted MDPs, given a dataset of i.i.d. trajectories. Under a similar linearity assumption, we provide a finite-sample OPE error bound for average-cost ergodic MDPs, and our approach only requires a single trajectory of the behavior policy.

**Maximum-entropy estimation.** The maximum-entropy principle has been well-studied in supervised learning (see e.g. Jaakkola et al. [2000]). There the objective is to maximize the entropy of a distribution subject to feature statistics matching on the available data, and the corresponding dual is maximum-likelihood estimation of an exponential family. In the batch RL setting, we maximize entropy subject to feature expectations matching under the MDP dynamics. For the linear MDP, the resulting distribution is also in the exponential family, and parameterized in a particular way that includes the model. Existing methods for modeling stationary distributions with function approximation tend to use linear functions and require extra constraints to ensure non-negativity and normalization [Rivera Cardoso et al., 2019, Abbasi-Yadkori et al., 2019b]. Exponential families seem like a more elegant solution, and also correspond to well-studied settings such as the linear quadratic regulator. Hazan et al. [2019] proposed learning maximum-entropy stationary distributions for the

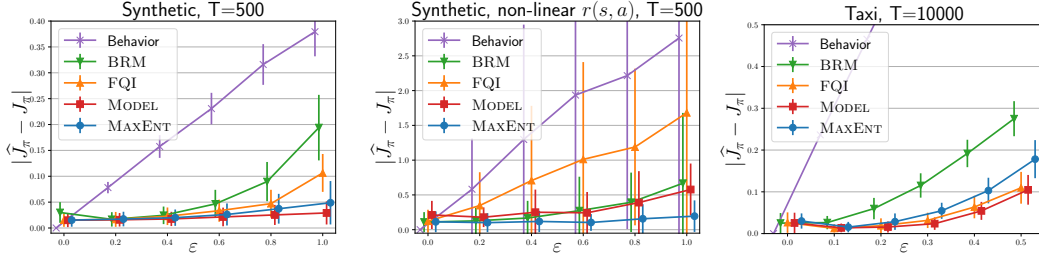


Figure 1: Experiments with behavior policy  $\varepsilon$ -greedy w.r.t.  $\pi$  on synthetic environments and Taxi (mean and standard deviation for 100 target policies  $\pi$ ). Note that the plots are slightly shifted along the horizontal axis to make error bars easier to see.

purpose of exploration. They focused on the tabular MDP case, and required an oracle for solving planning problems with function approximation, since in that case the entropy maximization problem may not be convex. We provide a convex formulation of this problem with function approximation in the linear MDP setting, which can also be used with neural networks (by learning representations).

## 5 Experiments

We compare our approach to other policy evaluation methods relying on function approximation. Since our focus is not on learning representations, we experiment with a fixed linear basis. We evaluate fitted Q-iteration (FQI) implemented as in (19) and Bellman residual minimization (BRM) implemented as in (18). We also use the average reward of the behavior policy as the simplest baseline. We refer to the closed-form version of our approach in (10) as MODEL, and to the version solving for the stationary distribution as MAXENT. We regularize the covariances of all regression problems using  $\alpha I$  with tuned  $\alpha$ .<sup>3</sup> For MAXENT, we optimize the parameters using full-batch Adam [Kingma and Ba, 2014], and normalize the distributions empirically. For experiments with OpenAI Gym environments [Brockman et al., 2016] (Taxi and Acrobot), we additionally use weighted importance sampling [Mahmood et al., 2014] for both the gradients and the objective. Unless stated otherwise, we generate policies by partially training on-policy using the POLITEX algorithm [Abbasi-Yadkori et al., 2019a], a version of regularized policy iteration with linear Q-functions. We compute the true policy values  $J_\pi$  using Monte-Carlo simulation for Acrobot, and exactly for other environments. Overall, we find that using a feature model is helpful with linear value-function methods.

**Synthetic environments.** We generate synthetic MDPs with 100 states, 10 actions, and transition matrices  $P$  generated by sampling entries uniformly at random and normalizing columns to sum to 1. We represent each state with a 10-dimensional vector  $\phi_S(s)$  of random Fourier features [Rahimi and Recht, 2008], and let  $\phi(s, a) = \phi_S(s) \otimes \phi_A(a)$ , where  $\phi_A(a)$  is a binary indicator vector for action  $a$ . We experiment with linear rewards  $r(s, a) = -\phi(s, a)^\top w$  with entries of  $w$  generated uniformly at random, and with non-linear rewards of the form  $r(s, a) = -\exp(2\phi(s, a)^\top w)$ . We generate target policies  $\pi$  by training on-policy, and set behavior policies  $\beta$  to be  $\varepsilon$ -greedy with respect to  $\pi$ . We plot the evaluation error  $|J_\pi - \hat{J}_\pi|$  for several values of  $\varepsilon$  in Figure 1 (showing mean and standard deviation for 100 random MDPs for each  $\varepsilon$ ). We can see that the model-based approaches are less sensitive to the difference between  $\pi$  and  $\beta$ , and the advantage of inferring the full distribution in the non-linear reward case. Note also that the true underlying dynamics are not low-rank, but our low-rank approximation still results in good estimates.

**Taxi.** The Taxi environment [Dietterich, 2000] is a  $5 \times 5$  grid with four pickup/dropoff locations. Taxi actions include going left, right, up, and down, and picking up or dropping off a passenger. There is a reward of -1 for every step, a reward of -10 for illegal pickup/dropoff actions, and a reward of 20 for a successful dropoff. In the infinite-horizon version, a new passenger appears after a successful dropoff. Our state features include indicators for whether the taxi is empty / at pickup / at dropoff and their pairwise products, and xy-coordinates of the taxi, passenger, and dropoff. We set  $\pi$  to be

<sup>3</sup>Starting with  $\alpha = 1$ , we keep doubling  $\alpha$  for FQI as long as it diverges, and for MAXENT as long as  $|\lambda_{\max}(\widehat{W}_\pi)| > 1$ .

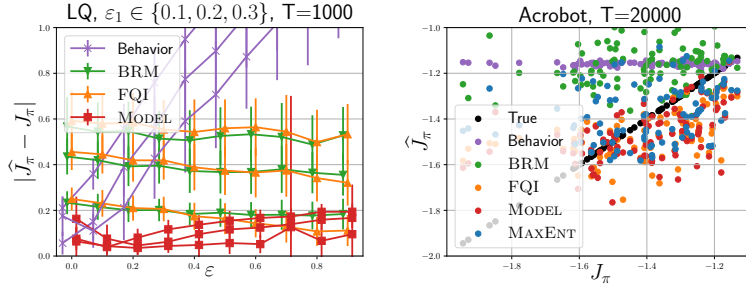


Figure 2: Left: experiments on LQ control, where  $\varepsilon$  and  $\varepsilon_1$  control the suboptimality of the behavior and target policies, respectively. The more suboptimal policies are difficult to evaluate with BRM and FQI. Right: predicted vs. true value on Acrobot for 100 target policies evaluated using the same behavior policy. The errors are: Behavior 0.24 ( $\pm 0.17$ ), BRM 0.23 ( $\pm 0.17$ ), FQI 0.14 ( $\pm 0.10$ ), MODEL 0.13 ( $\pm 0.11$ ), MAXENT 0.15 ( $\pm 0.14$ ).

0.05-greedy w.r.t. a hand-coded optimal strategy, and  $\beta$  to be  $\varepsilon$ -greedy w.r.t.  $\pi$ . A comparison of different policy evaluation methods is given in Figure 1. In this case, all methods are somewhat affected by the suboptimality of the behavior policy, possibly due to fewer successful dropoffs, and FQI and MODEL perform best.

**Linear quadratic regulator.** We evaluate our approach on the linear quadratic (LQ) control system in Dean et al. [2019], where stationary distributions, policies, and transition dynamics are Gaussian. We only evaluate a model-based approach here (as well as FQI and BRM) since the model fully constrains the solution, and solve a constrained optimization problem to ensure positive-definite covariances (see Appendix E for full details). We generate policies by solving the optimal control problem for the true dynamics and noisy costs, where  $\varepsilon$  controls the noise for  $\beta$  and  $\varepsilon_1$  controls the noise for  $\pi$ . The results are shown in Figure 2 (left) for ten values of  $\varepsilon$  and three values of  $\varepsilon_1$ . While  $\varepsilon$  does not seem to affect the OPE performance, the error increases with  $\varepsilon_1$  for BRM and FQI.

**Acrobot** [Sutton, 1996] is a simple episodic discrete-action physical control task. The system includes two links and two joints, one of which is actuated. The goal is to swing the lower link up to a given height. We set the reward at each time step to the negative distance between the joint to its target height, and to 100 when the lower link reaches its target height. Each episode ends after 500 steps, or after the target height is reached, after which we reset. The observations are link positions and velocities; we featurize them using the multivariate Fourier basis of order 3 as described in Konidaris et al. [2011]. For this task, we partially train 101 policies, set  $\beta$  to the first policy, and evaluate the remaining policies. The results are shown in Figure 2 (right). In this case, BRM predictions seem more correlated to the behavior policy than to the target. The other methods are better correlated with the target policy, but have somewhat high error. Possible reasons for this are the episodic nature of the environment, and the true underlying dynamics being only locally linear.

## 6 Conclusion and future work

We have presented a new approach to batch policy evaluation in average-reward MDPs. For linear MDPs, we have provided a finite-sample bound on the OPE error, which extends the previous results for discounted and episodic settings. In a more general setting with non-linear rewards and approximately linear feature dynamics, we have proposed a maximum-entropy approach to finding stationary distributions with function approximation. Given that the linear MDP assumption is fairly restrictive, one important direction for future work is extending the framework beyond linear functions. Another direction we are planning to explore is applying this framework to policy optimization. Finally, note that the maximum-entropy objective corresponds to minimizing the KL-divergence between the target distribution and the uniform distribution, and we can easily minimize the KL divergence to other distributions instead. While the maximum entropy objective is justified in some cases (see Appendix D), our formulation allows us to incorporate other prior knowledge and constraints when available, and this is another direction for future work.



## Broader impact

In general, when learning from a batch of data produced by a fixed behavior policy, we may inherit the biases of that policy, and our models may not generalize beyond the support of the data distribution. In our paper, we circumvent this issue by assuming that the information sufficient for evaluating and optimizing policies is contained in some known features, and that the behavior policy is exploratory enough in the sense that it spans those features. These assumptions may not always hold when applying the method in practice.

## References

- Yasin Abbasi-Yadkori, Peter Bartlett, Kush Bhatia, Nevena Lazić, Csaba Szepesvári, and Gellért Weisz. POLITEX: Regret bounds for policy iteration using expert prediction. In *Proceedings of the 36th International Conference on Machine Learning*, 2019a.
- Yasin Abbasi-Yadkori, Peter L. Bartlett, Xi Chen, and Alan Malek. Large-scale markov decision problems via the linear programming dual. *arXiv preprint arXiv:1901.01992*, 2019b.
- András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pages 30–37. Elsevier, 1995.
- Steven J Bradtke and Andrew G Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym, 2016.
- Jie Chen and Ronny Luss. Stochastic gradient descent with biased but consistent gradient estimators. *arXiv preprint arXiv:1807.11880*, 2018.
- Alon Cohen, Avinatan Hasidim, Tomer Koren, Nevena Lazić, Yishay Mansour, and Kunal Talwar. Online linear quadratic control. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1029–1038. PMLR, 10–15 Jul 2018.
- Christoph Dann, Gerhard Neumann, Jan Peters, et al. Policy evaluation with temporal differences: A survey and comparison. *Journal of Machine Learning Research*, 15:809–883, 2014.
- Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 2019.
- Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- Thomas G Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- Yaqi Duan and Mengdi Wang. Minimax-optimal off-policy evaluation with linear function approximation. *arXiv preprint arXiv:2002.09516*, 2020.
- Yihao Feng, Lihong Li, and Qiang Liu. A Kernel Loss for Solving the Bellman Equation. In *Advances in Neural Information Processing Systems 32*, pages 15456–15467. Curran Associates, Inc., 2019.
- Matthieu Geist and Bruno Scherrer. Off-policy learning with eligibility traces: A survey. *Journal of Machine Learning Research*, 15(1):289–333, 2014.

- Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy exploration. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2681–2691. PMLR, 2019.
- M Hazewinkel. Dirichlet distribution. *Encyclopedia of Mathematics*, 2001.
- Tommi Jaakkola, Marina Meila, and Tony Jebara. Maximum entropy discrimination. In *Advances in Neural Information Processing Systems*, pages 470–476, 2000.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- George Konidaris, Sarah Osentoski, and Philip Thomas. Value function approximation in reinforcement learning using the Fourier basis. In *25th AAAI Conference on Artificial Intelligence*, 2011.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pages 5356–5366, 2018.
- A Rupam Mahmood, Hado P van Hasselt, and Richard S Sutton. Weighted importance sampling for off-policy learning with linear function approximation. In *Advances in Neural Information Processing Systems*, pages 3014–3022, 2014.
- Ofir Nachum and Bo Dai. Reinforcement learning via fenchel-rockafellar duality. *arXiv preprint arXiv:2001.01866*, 2020.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections. In *Advances in Neural Information Processing Systems*, pages 2315–2325, 2019a.
- Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. AlgaeDICE: Policy Gradient from Arbitrary Experience. *arXiv preprint arXiv:1912.02074*, 2019b.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pages 1177–1184, 2008.
- Adrian Rivera Cardoso, He Wang, and Huan Xu. Large-scale Markov Decision Processes with changing rewards. In *Advances in Neural Information Processing Systems 32*, pages 2340–2350, 2019.
- RY Rubinstein. Simulation and the Monte Carlo method. 1981.
- Richard S Sutton. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In *Advances in Neural Information Processing Systems*, pages 1038–1044, 1996.
- Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th International Conference on Machine Learning*, pages 993–1000, 2009.
- Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- Junfeng Wen, Bo Dai, Lihong Li, and Dale Schuurmans. Batch stationary distribution estimation. *arXiv preprint arXiv:2003.00722*, 2020.
- Tengyang Xie and Nan Jiang. Q\* Approximation Schemes for Batch Reinforcement Learning: A Theoretical Comparison, 2020.

- Lin F Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019.
- Huizhen Yu. Convergence of Least Squares Temporal Difference Methods Under General Conditions. In *Proceedings of the International Conference on Machine Learning*, pages 1207–1214, 2010a.
- Huizhen Yu. Convergence of Least Squares Temporal Difference Methods Under General Conditions. Technical report, University of Helsinki, Department of Computer Science, 04 2010b.
- Huizhen Yu and Dimitri P Bertsekas. Convergence results for some temporal difference methods based on least squares. *IEEE Transactions on Automatic Control*, 54(7):1515–1531, 2009.
- Yi Yu, Tengyao Wang, and Richard J Samworth. A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika*, 102(2):315–323, 2015.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. GenDICE: Generalized offline estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020.

## A Off-policy evaluation dual objective

We formulate the estimation of the stationary state distribution  $\mu_\pi(s)$  given a policy  $\pi(a|s)$  as a maximum-entropy problem subject to matching feature expectations under the linearity assumption:

$$\min_{\mu} \sum_s \mu(s) \ln \mu(s) \quad (22)$$

$$\text{s.t. } \sum_{s,a} \mu(s) \pi(a|s) \phi(s, a)^\top (I - M_\pi) = b_\pi^\top \quad (23)$$

$$\sum_s \mu(s) = 1 \quad (24)$$

Let  $\phi(s, \pi) = \sum_a \pi(a|s) \phi(s, a)$  be the expected state features under the target policy. For a fixed  $I$  given  $M_\pi$ , the Lagrangian of the above objective is:

$$L(\mu, \theta, \lambda) = \sum_s \mu(s) \ln \mu(s) + b_\pi^\top \theta - \sum_s \mu(s) \phi(s, \pi)^\top (I - M_\pi) \theta - \lambda (\sum_s \mu(s) - 1)$$

Setting the gradient of  $L(\mu, \theta)$  w.r.t.  $\mu(s)$  to zero, we get

$$\begin{aligned} 0 &= \ln \mu(s) + 1 - \phi(s, \pi)^\top (I - M_\pi) \theta - \lambda \\ \mu(s) &= \exp(\phi(s, \pi)^\top (I - M_\pi) \theta + \lambda - 1) \end{aligned}$$

Because  $\sum_s \mu(s) = 1$ , we get

$$1 - \lambda = \ln \sum_s \exp(\phi(s, \pi)^\top (I - M_\pi) \theta) := F(\theta | M_\pi),$$

where  $F(\theta | M_\pi)$  is the log-normalizer. By plugging this expression for  $\mu$  into the Lagrangian, we get the following dual maximization objective in  $\theta$ :

$$\begin{aligned} D(\theta) &:= \sum_s \mu(s) (\phi(s, \pi)^\top (I - M_\pi) \theta - F(\theta | M_\pi)) + b_\pi^\top \theta - \sum_s \mu(s) \phi(s, \pi)^\top (I - M_\pi) \theta \\ &= b_\pi^\top \theta - F(\theta | M_\pi). \end{aligned}$$

## B Model error

### B.1 Preliminaries

Our error analysis relies on similar techniques as the finite-sample analysis in Abbasi-Yadkori et al. [2019a]. We first state some useful results.

**Lemma B.1** (Lemma A.1 in [Abbasi-Yadkori et al., 2019a]). *Let Assumption A1 hold, and let  $\{(s_t, a_t)\}_{t=1}^T$  be the state-action sequence obtained when following the behavior policy  $\beta$  from an initial distribution  $d_0$ . For  $t \in [T]$ , let  $X_t$  be a binary indicator vector with a non-zero at the linear index of the state-action pair  $(s_t, a_t)$ . Define for  $i \in [T]$ ,*

$$B_i = \mathbf{E} \left[ \sum_{t=1}^T X_t | X_1, \dots, X_i \right], \quad \text{and} \quad B_0 = \mathbf{E} \left[ \sum_{t=1}^T X_t \right].$$

*Then,  $(B_i)_{i=0}^T$  is a vector-valued martingale:  $\mathbf{E}[B_i - B_{i-1} | B_0, \dots, B_{i-1}] = 0$  for  $i = 1, \dots, T$ , and  $\|B_i - B_{i-1}\|_1 \leq 4\kappa$  holds for  $i \in [T]$ .*

The constructed martingale is known as the Doob martingale underlying the sum  $\sum_{t=1}^T X_t$ . Let  $\Pi_\beta$  be the transition matrix for state-action pairs when following  $\beta$ . Then, for  $t = 0, \dots, m-1$ ,  $\mathbf{E}[X_{t+1} | X_t] = \Pi_\beta^\top X_t$  and by the Markov property, for any  $i \in [T]$ ,

$$B_i = \sum_{t=1}^i X_t + \sum_{t=i+1}^T \mathbf{E}[X_t | X_i] = \sum_{t=1}^i X_t + \sum_{t=1}^{T-i} (\Pi_\beta^t)^\top X_i \quad \text{and} \quad B_0 = \sum_{t=1}^T (\Pi_\beta^t)^\top X_0.$$

It will be useful to define another Doob martingale as follows:

$$Y_i = \sum_{t=2}^i X_{t-1} X_t^\top + \sum_{t=i+1}^T \mathbf{E}[X_{t-1} X_t^\top | X_i] = \sum_{t=2}^i X_{t-1} X_t^\top + \sum_{t=i+1}^{T-i} \text{diag}(X_i^\top \Pi_\beta^{t-1}) \Pi_\beta \quad (25)$$

$$Y_0 = \sum_{t=1}^T \mathbf{E}[X_{t-1} X_t^\top] = \sum_{t=1}^T \text{diag}(d_0^\top \Pi_\beta^{t-1}) \Pi_\beta \quad (26)$$

where  $d_0$  is the initial state-action distribution. The difference sequence can again be bounded as  $\|Y_i - Y_{i-1}\|_{1,1} \leq 4\kappa$  under the mixing assumption (see Appendix D.2.2 of Abbasi-Yadkori et al. [2019a] for more details).

Let  $(\mathcal{F}_k)_k$  be a filtration and define  $\mathbf{E}_k[\cdot] := \mathbf{E}[\cdot | \mathcal{F}_k]$ .

**Theorem B.2** (Matrix Azuma [Tropp, 2012]). *Consider a finite  $(\mathcal{F})_k$ -adapted sequence  $\{X_k\}$  of Hermitian matrices of dimension  $m$ , and a fixed sequence  $\{A_k\}$  of Hermitian matrices that satisfy  $\mathbf{E}_{k-1} X_k = 0$  and  $X_k^2 \preceq A_k^2$  almost surely. Let  $v = \|\sum_k A_k^2\|$ . Then for all  $t \geq 0$ ,*

$$P\left(\lambda_{\max}\left(\sum_k X_k\right) \geq t\right) \leq m \cdot \exp(-t^2/8v).$$

Equivalently, with probability at least  $1 - \delta$ ,  $\|\sum_k X_k\| \leq 2\sqrt{2v \ln(m/\delta)}$ . A version of the inequality for non-Hermitian matrices of dimension  $m_1 \times m_2$  can be obtained by applying the theorem to a Hermitian dilation of  $X$ ,  $\mathcal{D}(X) = \begin{bmatrix} 0 & X \\ X^* & 0 \end{bmatrix}$ , which satisfies  $\lambda_{\max}(\mathcal{D}(X)) = \|X\|$  and  $\mathcal{D}(X)^2 = \begin{bmatrix} X X^* & 0 \\ 0 & X^* X \end{bmatrix}$ . In this case, we have that  $v = \max(\|\sum_k X_k X_k^*\|, \|\sum_k X_k^* X_k\|)$ .

Let  $\Phi$  be a  $|\mathcal{S}||\mathcal{A}| \times (m+1)$  matrix whose rows correspond to bias-augmented feature vectors  $\bar{\phi}(s, a)$  for each state-action pair  $(s, a)$ . Let  $\phi_i$  be the feature vector corresponding to the  $i^{\text{th}}$  row of  $\Phi$ , and let  $C_\Phi = \max_i \|\phi_i\|_2$ . For any matrix  $A$ , we have

$$\|\Phi^\top A \Phi\|_2 = \left\| \sum_{i,j} A_{ij} \phi_i \phi_j^\top \right\|_2 \leq \sum_{i,j} |A_{ij}| \|\phi_i \phi_j^\top\|_2 \leq C_\Phi^2 \sum_{i,j} |A_{ij}| = C_\Phi^2 \|A\|_{1,1}. \quad (27)$$

## B.2 Proof of Lemma 3.1

*Proof.* Let  $\Phi$  be a  $|\mathcal{S}||\mathcal{A}| \times (m+1)$  matrix whose rows correspond to bias-augmented feature vectors  $\bar{\phi}(s, a)$ . Let  $D_\beta = \text{diag}(d_\beta)$ . Let  $\tilde{d}_\beta$  be the empirical data distribution, and  $\tilde{D}_\beta = \text{diag}(\tilde{d}_\beta)$ . The true and estimated (concatenated) model parameters can be written as

$$\begin{aligned} \widehat{W}_\pi &= (\Lambda + \Phi^\top \tilde{D}_\beta \Phi)^{-1} \Phi^\top \tilde{D}_\beta \tilde{\Pi}_\pi \Phi \\ W_\pi &= (\Phi^\top D_\beta \Phi)^{-1} \Phi^\top D_\beta \Pi_\pi \Phi \end{aligned}$$

where  $\Pi_\pi$  is the true state-action transition kernel under  $\pi$ , and  $\tilde{\Pi}_\pi$  corresponds to empirical next-state dynamics  $\tilde{P}$ . For the true model satisfying  $\Pi_\pi \Phi = \Phi M_\pi$ , we have taken expectations over  $d_\beta$  and taken advantage of Assumption A3.

Let  $\Lambda = \alpha \Phi^\top \tilde{D}_\beta \Phi$ ; in this case

$$\widehat{W}_\pi = \frac{1}{1 + \alpha} (\Phi^\top \tilde{D}_\beta \Phi)^{-1} \Phi^\top \tilde{D}_\beta \tilde{\Pi}_\pi \Phi$$

We first bound the error for  $(1 + \alpha)\widehat{M}_\pi$ . The model error can be upper-bounded as:

$$\begin{aligned} \left\| (1 + \alpha)\widehat{W}_\pi - W_\pi \right\|_2 &\leq \left\| (\Phi^\top D_\beta \Phi)^{-1} \Phi^\top (\tilde{D}_\beta \tilde{\Pi}_\pi - D_\beta \Pi_\pi) \Phi \right\|_2 \\ &\quad + \left\| ((\Phi^\top \tilde{D}_\beta \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1}) \Phi^\top \tilde{D}_\beta \Pi_\pi \Phi \right\|_2 \\ &\leq \sigma^{-1} \left\| \Phi^\top (\tilde{D}_\beta - D_\beta) \Pi_\pi \Phi \right\|_2 + \left\| (\Phi^\top \tilde{D}_\beta \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1} \right\|_2 \left\| \Phi^\top \tilde{D}_\beta \Pi_\pi \Phi \right\|_2 \end{aligned}$$

$$\leq \sigma^{-1} \left\| \Phi^\top (\tilde{D}_\beta - D_\beta) \Pi_\pi \Phi \right\|_2 + C_\Phi^2 \left\| (\Phi^\top \tilde{D}_\beta \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1} \right\|_2$$

where the second inequality follows from Assumption A3, and the last inequality follows from (27). We proceed to bound the two terms

$$\begin{aligned} E_1 &= \sigma^{-1} \left\| \Phi^\top (\tilde{D}_\beta \tilde{\Pi}_\pi - D_\beta \Pi_\pi) \Phi \right\|_2 \\ E_2 &= \left\| (\Phi^\top \tilde{D}_\beta \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1} \right\|_2 \end{aligned}$$

**Bounding  $E_1$ .** Let  $(Y_i)_i$  be the Doob martingale defined in (25)-(26), and let  $\tilde{\Pi}_\beta$  be the empirical state-action transition matrix under the policy  $\beta$ . Note that  $\tilde{D}_\beta \tilde{\Pi}_\beta = Y_T/T$ . Furthermore, let  $K^\pi$  be a  $|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|$  matrix defined as

$$K_{(s,a),(s',a')}^\pi = \begin{cases} \pi(a'|s) & \text{if } s' = s \\ 0 & \text{otherwise} \end{cases}$$

Notice that  $\tilde{D}_\beta \tilde{\Pi}_\beta K^\pi = \tilde{D}_\beta \tilde{\Pi}_\pi$  and  $D_\beta \Pi_\beta K^\pi = D_\beta \Pi_\pi$ . We can upper-bound  $E_1$  as:

$$\sigma^{-1} \left\| \Phi^\top (\tilde{D}_\beta \tilde{\Pi}_\beta - D_\beta \Pi_\beta) K^\pi \Phi \right\|_2 = \frac{1}{\sigma T} \left\| \Phi^\top (Y_T - Y_0) K^\pi \Phi \right\|_2 + \frac{1}{\sigma T} \left\| \Phi^\top (Y_0 - T D_\beta \Pi_\beta) K^\pi \Phi \right\|_2$$

Note that  $\Phi^\top Y_i K^\pi \Phi$  is a matrix-valued martingale, whose difference sequence is bounded by

$$\left\| (\Phi^\top (Y_i - Y_{i-1}) K^\pi \Phi)^2 \right\|_2 \leq C_\Phi^4 \left\| (Y_i - Y_{i-1}) K^\pi \right\|_{1,1}^2 \leq 16 C_\Phi^4 \kappa^2$$

where we have used (27) and the fact that rows of  $K^\pi$  sum to 1. Applying the matrix-Azuma theorem B.2, we have that with probability at least  $\delta$ ,

$$\frac{1}{\sigma T} \left\| \Phi^\top (Y_T - Y_0) K^\pi \Phi \right\|_2 \leq 8 C_\Phi^2 \sigma^{-1} \kappa \sqrt{2 \ln(2(m+1)/\delta)/T}.$$

Using the mixing Assumption A1, and letting  $d_0$  be the initial state-action distribution,

$$\begin{aligned} \frac{1}{\sigma T} \left\| \Phi^\top (Y_0 - T D_\beta \Pi_\beta) K^\pi \Phi \right\|_2 &\leq \frac{1}{\sigma T} \sum_{t=1}^T \left\| \Phi^\top \text{diag}(d_0^\top \Pi_\beta^t - d_\beta^\top) \Pi_\beta K^\pi \Phi \right\|_2 \\ &\leq \frac{C_\Phi^2}{\sigma T} \sum_{t=1}^T \left\| \text{diag}(d_0^\top \Pi_\beta^t - d_\beta^\top) \Pi_\pi \right\|_{1,1} \\ &\leq \frac{C_\Phi^2}{\sigma T} \sum_{t=1}^T \exp(-t/\kappa) \|d_0 - d_\beta\|_1 \leq \frac{2 C_\Phi^2 \kappa}{\sigma T} \end{aligned}$$

Thus we get that with probability at least  $1 - \delta$ ,

$$E_1 \leq 8 C_\Phi^2 \sigma^{-1} \kappa \left( \sqrt{2 \ln(2(m+1)/\delta)/T} + 1/T \right)$$

**Bounding  $E_2$ .** To bound  $E_2$ , we first rely on the Woodbury identity to write

$$\begin{aligned} &(\Phi^\top \tilde{D}_\beta \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1} \\ &= (\Phi^\top D_\beta \Phi + \Phi^\top (D_\beta - \tilde{D}_\beta) \Phi)^{-1} - (\Phi^\top D_\beta \Phi)^{-1} \\ &= (\Phi^\top D_\beta \Phi)^{-1} \left( (\Phi^\top D_\beta \Phi)^{-1} + (\Phi^\top (\tilde{D}_\beta - D_\beta) \Phi)^{-1} \right) (\Phi^\top D_\beta \Phi)^{-1} \end{aligned}$$

$$\begin{aligned} E_2 &\leq \sigma^{-2} \left\| \left( (\Phi^\top D_\beta \Phi)^{-1} + (\Phi^\top (\tilde{D}_\beta - D_\beta) \Phi)^{-1} \right)^{-1} \right\|_2 \\ &\leq \sigma^{-2} \left\| \Phi^\top (\tilde{D}_\beta - D_\beta) \Phi \right\|_2 \\ &\leq 8 \sigma^{-2} C_\Phi^2 \kappa \left( \sqrt{2 \ln(2(m+1)/\delta)/T} + 1/T \right) \end{aligned}$$

where the second line follows because  $(\Phi^\top D_\beta \Phi)^{-1} \succ 0$ , and the last line follows by similar concentration arguments as those for  $E_1$  for the matrix-valued martingale  $\Phi^\top \text{diag}(B_i)\Phi$ , with probability at least  $1 - \delta$ .

**Bounding**  $\left\| \widehat{W}_\pi - W_\pi \right\|_2$ . Putting previous terms together, with probability at least  $1 - \delta$ , for an absolute constant  $C$ ,

$$\left\| (1 + \alpha) \widehat{W}_\pi - W_\pi \right\|_2 \leq CC_\Phi^4 \kappa \sigma^{-2} \sqrt{2 \ln(2(m+1)/\delta)/T} \quad (28)$$

Furthermore, we have:

$$\begin{aligned} \left\| \widehat{W}_\pi - W_\pi \right\|_2 &\leq \left\| (1 + \alpha) \widehat{W}_\pi - W_\pi \right\|_2 + \alpha \left\| \widehat{W}_\pi \right\|_2 \\ &\leq CC_\Phi^4 \kappa \sigma^{-2} \sqrt{2 \ln(2(m+1)/\delta)/T} + \alpha \sigma^{-1} C_\Phi^2 \end{aligned}$$

Setting  $\alpha = C_\Phi^2 \sigma^{-1} \kappa / \sqrt{T}$  gives the final result.

**Bounding**  $\|w - \hat{w}\|$ . For linear rewards  $r(s, a) = \phi(s, a)^\top w$ , we estimate the parameters  $w$  using linear regression. Abusing notation, assume that the feature matrix  $\Phi$  does not include bias for the purpose of this section. The true and estimated parameters  $w$  and  $\hat{w}$  satisfy

$$w = (\Phi^\top D_\beta \Phi)^{-1} \Phi^\top D_\beta r \quad (29)$$

$$\hat{w} = (\Phi^\top \tilde{D}_\beta \Phi)^{-1} \Phi^\top \tilde{D}_\beta r \quad (30)$$

where  $r = \Phi w$  is the length- $|\mathcal{S}||\mathcal{A}|$  vector of rewards. We have that

$$\|w - \hat{w}\| = \left\| (\Phi^\top D_\beta \Phi)^{-1} \Phi^\top (D_\beta - \tilde{D}_\beta) \Phi w \right\| + \left\| ((\Phi^\top D_\beta \Phi)^{-1} - (\Phi^\top \tilde{D}_\beta \Phi)^{-1}) \Phi^\top \tilde{D}_\beta \Phi w \right\| \quad (31)$$

Using the bounds from the previous section, we get that for a constant  $C_w$ , with probability at least  $1 - \delta$ ,

$$\|w - \hat{w}\| \leq C_w C_\Phi^2 \sigma^{-2} \kappa (\sqrt{2 \ln(2m/\delta)/T}) \|w\| \quad (32)$$

□

## C Proof of Theorem 3.2 (policy evaluation error)

*Proof.* Assuming that the optimization problem is feasible, the following holds for the resulting distribution  $\hat{d}_\pi(s, a) = \hat{\mu}_\pi(s) \pi(a|s)$ :

$$\sum_{s,a} \hat{d}_\pi(s, a) \bar{\phi}(s, a)^\top = \sum_{s,a} \hat{d}_\pi(s, a) \bar{\phi}(s, a)^\top \widehat{W}_\pi.$$

Assume that the reward is linear in the features,  $r(s, a) = w^\top \phi(s, a)$ , and let  $\hat{w}$  be the corresponding parameter estimate. Let  $\underline{w} = \begin{bmatrix} w \\ 0 \end{bmatrix}$  and let  $\underline{\hat{w}} = \begin{bmatrix} \hat{w} \\ 0 \end{bmatrix}$ .

The policy evaluation error is:

$$\begin{aligned} J_\pi - \hat{J}_\pi &= \sum_{s,a} (d_\pi(s, a) \bar{\phi}(s, a)^\top \underline{w} - \hat{d}_\pi(s, a) \bar{\phi}(s, a)^\top \underline{\hat{w}}) \\ &= \sum_{s,a} (d_\pi(s, a) - \hat{d}_\pi(s, a)) \bar{\phi}(s, a)^\top \underline{w} + \sum_{s,a} \hat{d}_\pi(s, a) \phi(s, a)^\top (\underline{w} - \underline{\hat{w}}). \end{aligned}$$

The norm of the second term is bounded by  $C_\Phi \|w - \hat{w}\|_2$ . We proceed to bound the first term.

Let  $W_\pi^\alpha = \frac{1}{1+\alpha} W_\pi$ , and define  $e^\top := \sum_{s,a} d(s, a) \bar{\phi}(s, a)^\top$  and  $\hat{e}^\top := \sum_{s,a} \hat{d}(s, a) \bar{\phi}(s, a)^\top$ . The first term can be written as:

$$(e^\top - \hat{e}^\top) \underline{w} = e^\top (W_\pi^\alpha + \alpha W_\pi^\alpha) w - \hat{e}^\top \widehat{W}_\pi \underline{w}$$

$$\begin{aligned}
&= (e - \hat{e})^\top W_\pi^\alpha \underline{w} + \hat{e}^\top (W_\pi^\alpha - \widehat{W}_\pi) \underline{w} + \alpha e^\top W_\pi^\alpha \underline{w} \\
&= (e - \hat{e})^\top (W_\pi^\alpha)^2 \underline{w} \\
&\quad + \hat{e}^\top (W_\pi^\alpha - \widehat{M}_\pi) (I + M_\pi) \underline{w} \\
&\quad + \alpha e^\top W_\pi^\alpha (I + (W_\pi^\alpha)^2) \underline{w} \\
&= \lim_{K \rightarrow \infty} (e - \hat{e})^\top (W_\pi^\alpha)^K \underline{w} + \left( \hat{e}^\top (W_\pi^\alpha - \widehat{W}_\pi) + \alpha e^\top W_\pi^\alpha \right) \left( \sum_{i=0}^K (W_\pi^\alpha)^i \right) \underline{w}
\end{aligned}$$

In order to evaluate the infinite sum, we first show that  $W_\pi$  is non-expansive in a  $\Sigma_\pi$ -weighted norm (and hence  $W_\pi^\alpha$  is contractive):

$$\begin{aligned}
\Sigma_\pi &:= \mathbf{E}_{(s,a) \sim d_\pi} [\overline{\phi}(s, a) \overline{\phi}(s, a)^\top] \\
&= \mathbf{E}_{(s,a) \sim d_\pi} [\mathbf{E}_{(s',a') \sim \Pi_\pi(\cdot|s,a)} [\overline{\phi}(s', a') \overline{\phi}(s', a')^\top]] \\
&= \mathbf{E}_{(s,a) \sim d_\pi} [W_\pi^\top \overline{\phi}(s, a) \overline{\phi}(s, a)^\top W_\pi] + V \\
&= W_\pi^\top \Sigma_\pi W_\pi + V
\end{aligned} \tag{33}$$

where  $V \succeq 0$ . Multiplying each side of (33) by  $\Sigma_\pi^{-1/2}$  from the left- and right-hand side, we get that

$$\begin{aligned}
I &= \Sigma_\pi^{-1/2} W_\pi^\top \Sigma_\pi^{1/2} \Sigma_\pi^{1/2} W_\pi \Sigma_\pi^{-1/2} + \Sigma_\pi^{-1/2} V \Sigma_\pi^{-1/2} \\
1 &\geq \left\| \Sigma_\pi^{1/2} W_\pi \Sigma_\pi^{-1/2} \right\|_2^2
\end{aligned}$$

Thus we have that  $\left\| \Sigma_\pi^{1/2} W_\pi^\alpha \Sigma_\pi^{-1/2} \right\|_2 \leq (1 + \alpha)^{-1}$ , and we can compute the infinite sum as:

$$\begin{aligned}
(W_\pi^\alpha)^i &= \Sigma_\pi^{-1/2} \left( \Sigma_\pi^{1/2} W_\pi^\alpha \Sigma_\pi^{-1/2} \right)^i \Sigma_\pi^{1/2} \\
\left\| \sum_{i=0}^{\infty} (W_\pi^\alpha)^i \right\| &\leq \sum_{i=0}^{\infty} \left\| \Sigma_\pi^{-1/2} \right\| \left\| \Sigma_\pi^{1/2} W_\pi^\alpha \Sigma_\pi^{-1/2} \right\|^i \left\| \Sigma_\pi^{1/2} \right\| \leq C_\Phi \sigma_\pi^{-1/2} (1 + \alpha)
\end{aligned}$$

The error can now be written as

$$|J_\pi - \widehat{J}_\pi| = \left( \|\hat{e}\|_2 \left\| W_\pi^\alpha - \widehat{W}_\pi \right\|_2 + \alpha \|e^\top W_\pi^\alpha\|_2 + \|u\|_2 \right) C_\Phi \sigma_\pi^{-1/2} (1 + \alpha) \|w\|_2 + C_\Phi \|w - \hat{w}\|_2$$

Note that  $\|e\|_2 \leq C_\Phi$  and  $\|\hat{e}\|_2 \leq C_\phi$ . Set  $\alpha = C_\Phi^2 \sigma_\pi^{-1} \kappa / \sqrt{T}$  as in the previous section. From (28), we have that

$$\left\| W_\pi^\alpha - \widehat{W}_\pi \right\|_2 \leq (1 + \alpha) C C_\Phi^4 \kappa \sigma_\pi^{-2} \sqrt{2 \ln(2(m+1)/\delta) / T} \tag{34}$$

Plugging in the model errors and  $\alpha$  and combining terms we get the final result in the theorem.  $\square$

## D Stationary distribution with large entropy

In this paper, we try to find the distribution over states that maximizes the entropy under some linear constraints, and use it as a proxy for the stationary distribution. In this section, we provide some theoretical evidence that at least when the probability transition over the states is sufficiently random, the stationary distribution tends to have large entropy.

For simplicity, we focus on finite-state Markov chains instead of MDPs. Consider a Markov chain with state space  $\mathcal{S}$  and probability transition matrix  $P$ . Let  $S := |\mathcal{S}|$ . Then the stationary distribution  $d$  satisfies  $d^\top = d^\top P$ . In this section, we assume that each row of  $P$  is sampled uniformly at random from the simplex over  $\mathcal{S}$ , i.e.,  $\Delta_S$ , independently of other rows. We prove the following result, which shows that as  $S$  increases, the stationary distribution  $d$  converges to a uniform distribution over the states at a rate  $\mathcal{O}(1/\sqrt{S})$ .



**Theorem D.1.** Let  $P$  be the probability transition matrix of a Markov chain with finite state space  $S$ , and assume that rows of  $P$  are sampled independently and uniformly at random from  $\Delta_S$ . Then, with probability at least  $1 - \delta$ , the stationary distribution  $d$  of the Markov chain satisfies

$$\left\| \frac{d}{\|d\|_2} - \frac{1}{\sqrt{S}} \mathbf{1} \right\|_2 \leq \frac{2\sqrt{10}}{\delta\sqrt{S}},$$

where  $\mathbf{1}$  denotes an all-one vector.

*Proof.* We denote the uniform distribution over the simplex in  $\mathbb{R}^S$  by  $\mathcal{U}$ . The distribution  $\mathcal{U}$  is a special case of Dirichlet distribution [Hazewinkel, 2001]. In this proof, we make use of the following properties of  $\mathcal{U}$ .

**Lemma D.2.** [Hazewinkel, 2001] Let  $x \sim \mathcal{U}$  and  $x_i$  be the  $i$ -th coordinate of  $x$ . Then we have

$$\begin{aligned} \mathbf{E}[x_i] &= \frac{1}{S}, & \mathbf{E}[x_i^2] &= \frac{2}{S(S+1)}, & \mathbf{E}[x_i^4] &= \frac{24}{S(S+1)(S+2)(S+3)} \\ \mathbf{E}[x_i x_j] &= \frac{1}{S(S+1)}, & \mathbf{E}[x_i^2 x_j^2] &= \frac{4}{S(S+1)(S+2)(S+3)}, & \forall i \neq j. \end{aligned}$$

This lemma gives us the following direct corollary.

**Corollary D.3.** Suppose that  $x$  and  $y$  are two independent samples from  $\mathcal{U}$ . Then we have

$$\mathbf{E}[\|x\|_2^2] = \frac{2}{S+1} \tag{35}$$

$$\mathbf{E}[x^\top y] = \frac{1}{S} \tag{36}$$

$$\mathbf{E}[\|x\|_2^4] = \frac{4(S+5)}{(S+1)(S+2)(S+3)} \tag{37}$$

$$\mathbf{E}[(x^\top y)^2] = \frac{S+3}{S(S+1)^2} \tag{38}$$

*Proof.*

$$\mathbf{E}[\|x\|_2^2] = \mathbf{E} \left[ \sum_{i=1}^S x_i^2 \right] = \frac{2}{S+1}.$$

$$\mathbf{E}[x^\top y] = \mathbf{E} \left[ \sum_{i=1}^S x_i y_i \right] = \sum_{i=1}^S \mathbf{E}[x_i] \mathbf{E}[y_i] = \frac{1}{S}.$$

$$\mathbf{E}[\|x\|_2^4] = \mathbf{E} \left[ \left( \sum_{i=1}^S x_i^2 \right)^2 \right] = \sum_{i=1}^S \mathbf{E}[x_i^4] + \sum_{i \neq j} \mathbf{E}[x_i^2 x_j^2] = \frac{4(S+5)}{(S+1)(S+2)(S+3)}.$$

$$\mathbf{E}[(x^\top y)^2] = \mathbf{E} \left[ \left( \sum_{i=1}^S x_i y_i \right)^2 \right] = \sum_{i=1}^S \mathbf{E}[x_i^2 y_i^2] + \sum_{i \neq j} \mathbf{E}[x_i x_j y_i y_j] = \frac{S+3}{S(S+1)^2}.$$

□

Now we turn to the proof of Theorem D.1. In the following, we define  $\widehat{\Sigma} := PP^\top$ ,  $\Sigma := \mathbf{E}[\widehat{\Sigma}]$ , and let  $p_i$  be the  $i$ -th column of  $P^\top$ . For a PSD matrix  $M$ , we define  $\lambda_i(M)$  as its  $i$ -th largest eigenvalue. Since  $P$  is a probability transition matrix, we know that  $\lambda_1(\widehat{\Sigma}) = 1$ , and the corresponding top eigenvector is  $\frac{d}{\|d\|_2}$ . We then analyze  $\Sigma$ . Since  $\Sigma_{i,j} = \mathbf{E}[p_i^\top p_j]$ , according to Corollary D.3, we know that  $\Sigma_{i,i} = \frac{2}{S+1}$ ,  $\forall i$  and  $\Sigma_{i,j} = \frac{1}{S}$ ,  $\forall i \neq j$ . Thus

$$\Sigma = \frac{S-1}{S(S+1)} I + \frac{1}{S} \mathbf{1}\mathbf{1}^\top.$$

Then, we know that  $\lambda_1(\Sigma) = 1 + \frac{S-1}{S(S+1)}$ ,  $\lambda_i(\Sigma) = \frac{S-1}{S(S+1)}$ ,  $\forall i \geq 2$ . Then, the gap between the top eigenvalue of  $\Sigma$  and the second largest eigenvalue of  $\Sigma$  is

$$\lambda_1(\Sigma) - \lambda_2(\Sigma) = 1. \quad (39)$$

The top eigenvector of  $\Sigma$  is  $\frac{1}{\sqrt{S}}\mathbf{1}$ . Next, we proceed to bound the difference between  $\Sigma$  and  $\widehat{\Sigma}$ . In particular, we bound  $\mathbf{E}[\|\widehat{\Sigma} - \Sigma\|_F^2]$ . We have

$$\begin{aligned} \mathbf{E}[\|\widehat{\Sigma} - \Sigma\|_F^2] &= \sum_{i=1}^S \left( \mathbf{E}[\widehat{\Sigma}_{i,i}^2] - \mathbf{E}[\widehat{\Sigma}_{i,i}]^2 \right) + \sum_{i \neq j} \left( \mathbf{E}[\widehat{\Sigma}_{i,j}^2] - \mathbf{E}[\widehat{\Sigma}_{i,j}]^2 \right) \\ &= \sum_{i=1}^S \left( \mathbf{E}[\|p_i\|_2^4] - \mathbf{E}[\|p_i\|_2^2]^2 \right) + \sum_{i \neq j} \left( \mathbf{E}[(p_i^\top p_j)^2] - \mathbf{E}[p_i^\top p_j]^2 \right) \\ &= \frac{4S(S-1)}{(S+1)^2(S+2)(S+3)} + \frac{(S-1)^2}{S(S+1)^2} \end{aligned} \quad (40)$$

$$\leq \frac{5}{S}, \quad (41)$$

where in (40) we use Corollary D.3. Thus, we have

$$\mathbf{E}[\|\widehat{\Sigma} - \Sigma\|_F] \leq \sqrt{\mathbf{E}[\|\widehat{\Sigma} - \Sigma\|_F^2]} \leq \sqrt{\frac{5}{S}}. \quad (42)$$

According to Markov's inequality, with probability at least  $1 - \delta$ ,

$$\|\widehat{\Sigma} - \Sigma\|_F \leq \frac{\sqrt{5}}{\delta\sqrt{S}}. \quad (43)$$

We then apply Davis-Kahan Theorem [Davis and Kahan, 1970] (see also Theorem 2 in Yu et al. [2015]) and obtain

$$\sqrt{1 - \left\langle \frac{d}{\|d\|_2}, \frac{1}{\sqrt{S}}\mathbf{1} \right\rangle^2} \leq \frac{2\|\widehat{\Sigma} - \Sigma\|_F}{\lambda_1(\Sigma) - \lambda_2(\Sigma)} = 2\|\widehat{\Sigma} - \Sigma\|_F,$$

where for the equality we use (39). This implies

$$\begin{aligned} \left\| \frac{d}{\|d\|_2} - \frac{1}{\sqrt{S}}\mathbf{1} \right\|_2 &= \sqrt{2 - 2\left\langle \frac{d}{\|d\|_2}, \frac{1}{\sqrt{S}}\mathbf{1} \right\rangle} \\ &\leq \sqrt{2} \sqrt{1 - \left\langle \frac{d}{\|d\|_2}, \frac{1}{\sqrt{S}}\mathbf{1} \right\rangle^2} \\ &\leq 2\sqrt{2}\|\widehat{\Sigma} - \Sigma\|_F. \end{aligned} \quad (44)$$

Then we can complete the proof by combining (43) and (44).  $\square$

## E Experiment details for linear quadratic control

In a linear-quadratic (LQ) control problem, the dynamics are linear-Gaussian in states  $x$ :

$$x_{t+1} = Ax_t + Ba_t + w_t, \quad w_t \sim \mathcal{N}(0, W). \quad (45)$$

Assume that all policies are linear-Gaussian:  $\pi(a|x) = \mathcal{N}(a|Kx, C)$ . In this case, assuming that the policy  $\pi$  is stable (the spectral radius of  $A + BK$  is less than 1), the stationary state distribution is

$$\mu(x) = \mathcal{N}(0, S), \quad \text{where } S = (A + BK)S(A + BK)^\top + W. \quad (46)$$

Given an estimate of the dynamics parameters  $(\widehat{A}, \widehat{B}, \widehat{W})$ , maximum-entropy OPE corresponds to the following convex problem:

$$\max_{S \succeq 0} \ln \det(S) \quad (47)$$

$$\text{s.t. } S = (\widehat{A} + \widehat{B}K)S(\widehat{A} + \widehat{B}K)^\top + \widehat{W} \quad (48)$$

Note that the constraint corresponds to that in the dual formulation of LQ control presented in Cohen et al. [2018]. We solve the above problem using cvxpy [Diamond and Boyd, 2016]. The problem will only be feasible if  $\rho(\widehat{A} + \widehat{B}K) < 1$ , where  $\rho(\cdot)$  denotes the spectral radius. Furthermore, when the system is controllable, the constraint fully specifies the solution and so the maximum-entropy objective plays no role.

In LQ control problems, rewards are quadratic:

$$r(x, a) = -x^\top Qx - a^\top Ra, \quad Q, R \succ 0, \quad (49)$$

Thus to evaluate policies, we can estimate  $Q$  and  $R$ , and estimate the policy value as

$$\widehat{J}_\pi = \text{trace}(S\widehat{Q}) + \text{trace}((KSK^\top + C)\widehat{R}).$$

In our experimental setup, we produce the behavior policies by solving for the optimal controller for true dynamics  $(A, B, W)$ , true action costs  $R$ , and state costs corrupted as

$$\widetilde{Q} = Q + \varepsilon^2 U^\top U$$

$U$  is a matrix of the same size as  $Q$  whose entries are generated uniformly at random. Given the corresponding optimal linear feedback matrices  $\widetilde{K}$ , we set behavior policies to  $\beta(a|x) = \mathcal{N}(a|\widetilde{K}x, 0.1I)$ , and we make target policies greedy, i.e.  $a = \widetilde{K}x$ .

When evaluating policies using BRM and FQI, we use the following features for a policy  $\pi(a|x) = \mathcal{N}(a|Kx, C)$ :

$$\phi(s, a) = \text{VEC} \left( \begin{bmatrix} xx^\top & xa^\top \\ ax^\top & aa^\top \end{bmatrix} \right), \quad \phi(s, \pi) = \text{VEC} \left( \begin{bmatrix} xx^\top & xx^\top K^\top \\ Kxx^\top & Kxx^\top K^\top + C \end{bmatrix} \right).$$

## F Batch policy optimization

To optimize policies, we can maximize the entropy-regularized expected reward  $\sum_{s,a} d(s, a)(r(s, a) - \tau \ln d(s, a))$  subject to the same feature constraints as before. Unfortunately, in this case the feature expectation constraints are no longer linear, as the model  $\widehat{M}_\pi$  depends on the optimization variables through  $\pi(a|s) = \frac{d(s, a)}{\sum_{a'} d(s, a')}$ . One possible optimization approach is an EM-like algorithm that alternates between optimizing  $d(s, a)$  for fixed  $\widehat{M}_\pi$ , and reestimating  $\widehat{M}_\pi$  for  $\pi(a|s) \propto d(s, a)$ . A simpler alternative, proposed in Yang and Wang [2019], is to assume that we have state-only features  $\psi(s)$  whose expectation is a linear function of the state-action features:

$$\mathbf{E}_{s' \sim P(\cdot|s, a)}[\psi(s')] = \phi(s, a)^\top M$$

where  $M$  is a matrix of appropriate dimensions. Note that now  $M$  does not depend on the policy, and can be kernelized as in Yang and Wang [2019]. With this, we formulate batch policy optimization as:

$$\min_{d \in \Delta_{\mathcal{S} \times \mathcal{A}}} \sum_{s, a} d(s, a)(-\phi(s, a)^\top \hat{w} + \tau \ln d(s, a)) \quad (50)$$

$$\text{s.t. } \sum_{s, a} d(s, a)\phi(s, a)^\top \widehat{M} = \sum_{s, a} d(s, a)\psi(s)^\top \quad (51)$$

The optimal solution takes the form

$$d(s, a|\theta_d, \hat{w}, \widehat{M}) = \exp \left( \frac{1}{\tau} \phi(s, a)^\top \hat{w} + \frac{1}{\tau} (\phi(s, a)^\top \widehat{M} - \psi(s)^\top) \theta_d - F_\tau(\theta_d, \hat{w}, \widehat{M}) \right) \quad (52)$$

where  $F_\tau(\theta_d, \hat{w}, \widehat{M})$  is the log-partition function, and  $\theta_d = \arg \min_\theta F(\theta, \hat{w}, \widehat{M})$ .