
Learning Bellman Complete Representations for Offline Policy Evaluation

Jonathan D. Chang^{*1} Kaiwen Wang^{*1} Nathan Kallus² Wen Sun¹

Abstract

We study representation learning for Offline Reinforcement Learning (RL), focusing on the important task of Offline Policy Evaluation (OPE). Recent work shows that, in contrast to supervised learning, realizability of the Q-function is not enough for learning it. Two sufficient conditions for sample-efficient OPE are Bellman completeness and coverage. Prior work often assumes that representations satisfying these conditions are given, with results being mostly theoretical in nature. In this work, we propose BCRL, which directly learns from data an approximately linear Bellman complete representation with good coverage. With this learned representation, we perform OPE using Least Square Policy Evaluation (LSPE) with linear functions in our learned representation. We present an end-to-end theoretical analysis, showing that our two-stage algorithm enjoys polynomial sample complexity provided *some* representation in the rich class considered is linear Bellman complete. Empirically, we extensively evaluate our algorithm on challenging, image-based continuous control tasks from the Deepmind Control Suite. We show our representation enables better OPE compared to previous representation learning methods developed for off-policy RL (e.g., CURL, SPR). BCRL achieves competitive OPE error with the state-of-the-art method Fitted Q-Evaluation (FQE), and beats FQE when evaluating beyond the initial state distribution. Our ablations show that both linear Bellman complete and coverage components of our method are crucial.

1. Introduction

Deep Reinforcement Learning (RL) has developed agents that solve complex sequential decision making tasks, achieving new state-of-the-art results and surpassing expert human performance. Despite these impressive results, these algorithms often require a prohibitively large number of online interactions to scale to higher dimensional inputs.

To address these sample complexity demands, a recent line of work (van den Oord et al., 2018; Anand et al., 2019; Mazouze et al., 2020; Yang & Nachum, 2021) has incorporated advances in unsupervised representation learning from the supervised learning literature into developing RL agents. For example, CURL (Laskin et al., 2020) and SPR (Schwarzer et al., 2020) utilize contrastive representation objectives as auxiliary losses within an existing RL framework. These efforts are motivated by the tremendous success such self-supervised techniques have offered in computer vision, natural language processing, speech processing, and beyond. While these formulations have shown sample complexity improvements empirically, it remains an open question whether these approaches successfully address the unique challenges from RL that usually do not appear in supervised learning, such as exploration and exploitation, credit assignments, long horizon prediction, and distribution shift. In particular, recent work (Wang et al., 2021; Amortila et al., 2020; Foster et al., 2021) has shown that realizability of the learning target in RL (namely, the Q-function) is insufficient to avoid exponential dependence on problem parameters.

In this paper, we study representation learning for an important subtask of off-policy RL: offline policy evaluation (OPE). OPE is a critical component for any off-policy policy optimization approach (e.g., off-policy actor critic such as SAC, Haarnoja et al., 2018, and DDPG, Lillicrap et al., 2016). Moreover, OPE allows us to focus on issues arising from distribution shift and long horizon prediction. Specifically, we propose a new approach that leverages rich function approximation (e.g., deep neural networks) to learn a representation that is both Bellman complete and exploratory. A linear Bellman complete representation means that linear functions in the representation have zero inherent Bellman error (Antos et al., 2008), i.e., applying the Bellman operator on a linear function results in a new linear function. An exploratory representation means that

^{*}Equal contribution ¹Computer Science, Cornell University, Ithaca, NY, USA ²Operations Research and Information Engineering, Cornell Tech, New York, NY, USA. Correspondence to: Jonathan D. Chang <<https://jdchang1.github.io>>, Kaiwen Wang <<https://kaiwenw.github.io>>.

the resulting feature covariance matrix formed by the offline dataset is well-conditioned. These two representational properties ensure that, under linear function approximation (i.e., the linear evaluation protocol, Grill et al., 2020), classic least squares policy evaluation (LSPE) (Nedić & Bertsekas, 2003; Duan et al., 2020; Wang et al., 2020) can achieve accurate policy evaluation. We provide an end-to-end analysis showing that our representation learning approach together with LSPE ensures near-optimal policy evaluation with polynomial sample complexity.

Empirically, we extensively evaluate our method on image-based continuous control tasks from the Deepmind Control Suite. First, we compare against two representation learning approaches developed for off-policy RL: CURL (Laskin et al., 2020) and SPR (Schwarzer et al., 2020). These bear many similarities to contrastive learning techniques for unsupervised representation learning: SimCLR (Chen et al., 2020) and Bootstrap your own latent (BYOL, Grill et al., 2020), respectively. Under the linear evaluation protocol (i.e., LSPE with a linear function on top of the learned representation), our approach consistently outperforms these baselines. We observe that representations learned by CURL and SPR sometimes even exhibit *instability when evaluated using LSPE* (prediction error blows up when more iterations of LSPE is applied), while our approach is more stable. This comparison demonstrates that representation learning in offline RL is more subtle and using representation techniques developed from supervised learning settings may not result in the best performance for offline RL. Our ablations show that both linear Bellman completeness and coverage are crucial, as our method also blows up if one ingredient is missing. Finally, BCRL achieves state-of-the-art OPE error when compared with other OPE methods, and improves the state-of-the-art when evaluating beyond the initial state distribution.¹

1.1. Related Work

Representation Learning in Offline RL:

From the theoretical side, Hao et al. (2021) considers offline RL in sparse linear MDPs (Jin et al., 2020). Learning with sparsity can be understood as feature selection. Their work has a much stronger coverage condition than ours: namely, given a representation class \mathcal{C} , they assume that any feature $\phi \in \mathcal{C}$ has global coverage under the offline data distribution, i.e., $\mathbb{E}_{s,a} \nu \phi(s,a) \phi(s,a)^\top$ is well conditioned where ν is offline data distribution. In our work, we only assume that *there exists one* ϕ^* that has global coverage, thus strictly generalizing their coverage condition. Uehara & Sun (2021) propose a general model-based offline RL approach that can perform representation learning for linear MDPs in the offline setting without global coverage. How-

ever, their algorithm is a version space approach and is not computationally efficient. Also, our linear Bellman completeness condition *strictly generalizes* linear MDPs. (Ni et al., 2021) consider learning state-action embeddings from a known RKHS. We use general function approximation that can be more powerful than RKHS. Finally, Parr et al. (2008) identifies bellman completeness as a desirable condition for feature selection in RL when analyzing an equivalence between linear value-function approximation and linear model approximation. In our work, we investigate how to learn bellman completeness representation, and also the role of coverage in our feature selection.

From the empirical side, Yang & Nachum (2021) evaluated a broad range of unsupervised objectives for learning pretrained representations from offline datasets for downstream Imitation Learning (IL), online RL, and offline RL. They found that the use of pretrained representations dramatically improved the downstream performance for policy learning algorithms. In this work, we aim to identify such an unsupervised representation objective and to extend the empirical evaluation to offline image-based continuous control tasks. Nachum & Yang (2021) presents a provable contrastive representation learning objective for IL (derived from maximum likelihood loss), learning state representations from expert data to do imitation with behavior cloning. Our approach instead focuses on learning state-action representations for OPE. Finally, both Song et al. (2016) and Chung et al. (2019) present algorithms to learn representations suitable for linear value function approximation. Song et al. (2016) is most relevant to our work where they aim to learn bellman complete features. Both works, however, work in the online setting and do not consider coverage induced from the representation which we identify is important for accurate OPE.

Representation Learning in Online RL:

Theoretical works on representation learning in online RL focus on learning representations to facilitate exploration from scratch. OLIVE (Jiang et al., 2017), BLin-UCB (Du et al., 2021), FLAMBE (Agarwal et al., 2020), Moflle (Modi et al., 2021), and Rep-UCB (Uehara et al., 2022) propose approaches for representation learning in linear MDPs where they assume that there exists a feature $\phi^* \in \mathcal{C}$ that admits the linear MDP structure for the ground truth transition. Zhang et al. (2021) posits an even stronger assumption where every feature $\phi \in \mathcal{C}$ admits the linear MDP structure for the true transition. Note that we only assume that *there exists one* ϕ^* that admits linear Bellman completeness, which strictly generalizes linear MDPs. Hence, our representation learning setting is more general than prior theoretical works. Note that we study the offline setting while prior works operate in the online setting which has additional challenges from online exploration.

¹Code available at <https://github.com/CausalML/bcrl>.

On the empirical side, there is a large body of works that adapt existing self-supervised learning techniques developed from computer vision and NLP to RL. For learning representations from image-based RL tasks, CPC (van den Oord et al., 2018), ST-DIM (Anand et al., 2019), DRIML (Mazouze et al., 2020), CURL (Laskin et al., 2020), and SPR (Schwarzer et al., 2020) learn representations by optimizing various temporal contrastive losses. In particular, Laskin et al. (2020) proposes to interleave minimizing an InfoNCE objective with similarities to MoCo (He et al., 2019) and SimCLR (Chen et al., 2020) while learning a policy with SAC (Haarnoja et al., 2018). Moreover, Schwarzer et al. (2020) proposes learning a representation similar to BYOL (Grill et al., 2020) alongside a Q-function for Deep Q-Learning (Mnih et al., 2013). We compare our representational objective against CURL and SPR in Section 6, and demonstrate that under linear evaluation protocol, ours outperform both CURL and SPR. Note, we refer the readers to Schwarzer et al. (2020) for comparisons between SPR and CPC, ST-DIM, and DRIML.

2. Preliminaries

In this paper we consider the infinite horizon discounted MDP $(S, A, \gamma, P, r, d_0)$ where S, A are state and action spaces which could contain a large number of states and actions or could even be infinite, $\gamma \in (0, 1)$ is a discount factor, P is the transition kernel, $r : S \times A \rightarrow \mathbb{R}$ is the reward function, and $d_0 \in \Delta(S)$ is the initial state distribution. We assume rewards are bounded by 1, i.e. $|r(s, a)| \leq 1$. Given a policy $\pi : S \times A \rightarrow \Delta(A)$, we denote $V^\pi(s) = \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h r_h | \pi, s_0 := s]$ as the expected discounted total reward of policy π starting at state s . We denote $V_{d_0}^\pi = \mathbb{E}_{s \sim d_0} V^\pi(s)$ as the expected discounted total reward of the policy π starting at the initial state distribution d_0 . We also denote average state-action distribution $d_{d_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h d_h^\pi(s, a)$ where $d_h^\pi(s, a)$ is the probability of π visiting the (s, a) pair at time step h , starting from d_0 .

In OPE, we seek to evaluate the expected discounted total reward V^{π_e} of a target policy $\pi_e : S \times A \rightarrow \Delta(A)$ given data drawn from an offline state-action distribution $\nu \in \Delta(S \times A)$. The latter can, for example, be the state-action distribution under a behavior policy π_b . Namely, the offline dataset $D = \{s_i, a_i, r_i, s_i^0\}_{i=1}^N$ consists of N i.i.d tuples generated as $(s, a) \sim \nu, r = r(s, a), s^0 \sim P(\cdot | s, a)$.

We define Bellman operator T^π associated with π as follows:

$$T^\pi f(s, a) = r(s, a) + \gamma \mathbb{E}_{s^0 \sim P(s, a), a^0 \sim \pi(s^0)} [f(s^0, a^0)]$$

We may drop the superscript π when it is clear from context, in particular when $\pi = \pi_e$ is the fixed target policy.

A representation, or feature, $\phi : S \times A \rightarrow \mathbb{R}^d$ is an embedding of state-action pairs into d -dimensional

space. We let $\|\phi\| = \mathbb{E}_\nu [\|\phi(s, a)\|]$, $\widehat{\|\phi\|} = \frac{1}{N} \sum_{i=1}^N \|\phi(s_i, a_i)\|$. We consider learning a representation from a feature hypothesis class $\mathcal{F} = \{S \times A \rightarrow \mathbb{R}^d\}$. We assume features have bounded norm: $\|\phi(s, a)\| \leq 1$. In our experiments, ϕ is convolutional neural nets with d outputs.

Notation We denote $B_W := \{x \in \mathbb{R}^d : \|x\| \leq W\}$ as the Euclidean Ball in \mathbb{R}^d with radius W . Given a distribution $\nu \in \Delta(S \times A)$ and a function $f : S \times A \rightarrow \mathbb{R}$, we denote $L_2(\nu)$ norm as $\|f\|_{L_2(\nu)} = \sqrt{\mathbb{E}_{s, a \sim \nu} f^2(s, a)}$. Given a positive definite matrix Σ , let $\|x\|_\Sigma = \sqrt{x^T \Sigma x}$ and let $\lambda_{\min}(\Sigma)$ denote the minimum eigenvalue. When $\nu \ll \mu$ we let $\frac{d\nu}{d\mu}$ denote the Radon-Nikodym derivative. We use \circ to denote composition, so $s^0, a^0 \sim P(s, a) \circ \pi$ is short-form for $s^0 \sim P(s, a), a^0 \sim \pi(s^0)$. For any function $f(s, a)$ and a policy π , we denote $f(s, \pi) = \mathbb{E}_a \pi(s) [f(s, a)]$.

3. Linear Bellman Completeness

Before we introduce our representation learning approach, in this section, we first consider OPE with linear function approximation with a given representation ϕ . Lessons from supervised learning or linear bandit may suggest that OPE should be possible with polynomially many offline samples, as long as (1) ground truth target Q^{π_e} is linear in ϕ (i.e. $Q^{\pi_e} \in \mathbb{R}^d$, such that for all $s, a, Q^{\pi_e}(s, a) = w^T \phi(s, a)$), and (2) the offline data provides sufficient coverage over π_e (i.e., $\lambda_{\min}(\mathbb{E}_{\nu} [\phi \phi^T]) > 0$). Unfortunately, under these two assumptions, there are lower bounds indicating that for any OPE algorithm, there exists an MDP where one will need at least exponentially (exponential in horizon or d) many offline samples to provide an accurate estimate of V^{π_e} (Wang et al., 2021; Foster et al., 2021; Amortila et al., 2020). This lower bound indicates that one needs additional structural conditions on the representation.

The additional condition the prior work has considered is Bellman completeness (BC). Since we seek to learn a representation rather than assume theoretical conditions on a given one, we will focus on an *approximate* version of BC.

Definition 3.1 (Approximate Linear BC). A representation ϕ is ε_ν -approximately linear Bellman complete if,

$$\max_{w_1 \in B_W} \min_{w_2 \in B_W} \|w_1 - T^{\pi_e}(w_2)\|_{L_2(\nu)} \leq \varepsilon_\nu.$$

Note the dependence on ν, π_e , and W .

Intuitively the above condition is saying that for any linear function $w_1^T \phi(s, a)$, $T^\pi(w_1^T \phi(s, a))$ itself can be approximated by another linear function under the distribution ν .

Remark 3.2. Low-rank MDPs (Jiang et al., 2017; Agarwal et al., 2020) are subsumed in the exact linear BC model, i.e., $\varepsilon_\nu = 0$ under any distribution ν .

Algorithm 1 Least Squares Policy Evaluation (LSPE)

- 1: **Input:** Target policy π_e , features ϕ , dataset D
- 2: Initialize $\hat{\theta}_0 = \mathbf{0} \succeq B_W$.
- 3: **for** $k = 1, 2, \dots, K$ **do**
- 4: Set $\hat{f}_{k-1}(s, a) = \hat{\theta}_{k-1}^\top \phi(s, a)$,
 $\hat{V}_{k-1}(s) = \mathbb{E}_{a \sim \pi_e(s)} [\hat{f}_{k-1}(s, a)]$
- 5: Perform linear regression:

$$\hat{\theta}_k \succeq \arg \min_{\theta \succeq B_W} \frac{1}{N} \sum_{i=1}^N (\theta^\top \phi(s_i, a_i) - r_i - \gamma \hat{V}_{k-1}(s_i^\theta))^2$$
- 6: **end for**
- 7: Return \hat{f}_K .

Note that the Bellman completeness condition is more subtle than the common realizability condition: for any fixed ϕ , increasing its expressiveness (e.g., add more features) does not imply a monotonic decrease in ε_ν . Thus common tricks such as lifting features to higher order polynomial kernel space or more general reproducing kernel Hilbert space does not imply the linear Bellman complete condition, nor does it improve the coverage condition. We next show approximate Linear BC together with coverage imply sample-efficient OPE via Least Square Policy Evaluation (LSPE) (Algorithm 1). We present our result using the relative condition number, defined for any initial state distribution p_0 as

$$\kappa(p_0) := \sup_{x \in \mathbb{R}^d} \frac{x^\top \mathbb{E}_{s,a} \frac{d \pi_e}{d p_0} \phi(s, a) \phi(s, a)^\top x}{x^\top (\phi) x}. \quad (1)$$

Theorem 3.3 (Sample Complexity of LSPE). *Assume feature ϕ satisfies approximate Linear BC with parameter ε_ν . For any $\delta \succeq (0, 0.1)$, w.p. at least $1 - \delta$, we have for any initial state distribution p_0*

$$\left\| V_{p_0}^{\pi_e} - \mathbb{E}_{s \sim p_0} [\hat{f}_K(s, \pi_e)] \right\| \leq \frac{\gamma^{K/2}}{1 - \gamma} + \frac{4 \sqrt{\left\| \frac{d \pi_e}{d p_0} \right\|_1}}{(1 - \gamma)^2} \varepsilon_\nu + \frac{480 \sqrt{\kappa(p_0)} (1 + W) d \log(N) \sqrt{\log(10/\delta)}}{(1 - \gamma)^2 \beta \overline{N}},$$

where \hat{f}_K is the output of Algorithm 1.

Please see Appendix C.1 for proof. The above result holds simultaneously over all initial state distributions p_0 covered by the data distribution ν . If ν has full coverage, i.e. if βI , as is commonly assumed in the literature, one can show that $\kappa(p_0) \leq \beta^{-1}$ for any initial state distribution p_0 . Also note that the concentrability coefficient $\left\| \frac{d \pi_e}{d \nu} \right\|_1$ shows up as $T^{\pi_e} \hat{f}_{k-1}$ can be *nonlinear*. In the exact Linear BC case, where $\varepsilon_\nu = 0$ (i.e., there is a linear function that perfectly approximates $T^{\pi_e} \hat{f}_{k-1}$ under ν), the term involving the concentrability coefficient will be 0 and we can even avoid its finiteness.

4. Representation Learning

The previous section indicates sufficient conditions on the representation for efficient and accurate OPE with linear function approximation. However, a representation that is Bellman complete and also provides coverage is not available a priori, especially for high-dimensional settings (e.g., image-based control tasks as we consider in our experiments). Existing theoretical work often assume such representation is given. We propose to learn a representation ϕ that is approximately Bellman complete and also provides good coverage, via rich function approximation (e.g., a deep neural network). More formally, we want to learn a representation ϕ such that (1) it ensures approximate Bellman complete, i.e., ε_ν is small, and (2) has a good coverage, i.e., $\lambda_{\min}(\phi)$ is as large as possible, which are the two key properties to guarantee accurate and efficient OPE indicated by Theorem 3.3. To formulate the representation learning question, our key assumption is that our representation hypothesis class is rich enough such that it contains at least one representation ϕ^* that is linear Bellman complete (i.e., $\varepsilon_\nu = 0$) and has a good coverage.

Assumption 4.1 (Representational power of \mathcal{H}). There exists $\phi^* \in \mathcal{H}$, such that (1) ϕ^* achieves exact Linear BC (definition 3.1 with $\varepsilon_\nu = 0$), and (2) ϕ^* induces coverage, i.e., $\lambda_{\min}(\phi^*) \geq \beta > 0$.

Our goal is to learn such a representation from the hypothesis class \mathcal{H} . Note that unlike prior RL representation learning works (Hao et al., 2021; Zhang et al., 2021), here we only assume that there exists a ϕ^* has linear BC and induces good coverage, other candidate $\phi \in \mathcal{H}$ could have terrible coverage and does not necessarily have linear BC.

Before proceeding to the learning algorithm, we first present an *equivalent condition* for linear BC, which does not rely on the complicated min-max expression of Definition 3.1.

Proposition 4.2. *Consider a feature ϕ with full rank covariance $\Sigma(\phi)$. Given any $W > 0$, the feature ϕ being linear BC (under B_W) implies that there exist $(\rho, M) \succeq B_W \in \mathbb{R}^d \times \mathbb{R}^d$ with $k M k_2 \leq \sqrt{1 - \frac{k \rho k_2^2}{W^2}}$, and*

$$\mathbb{E}_\nu \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \mathbb{E}_{s^\theta} P(s, a) \phi(s^\theta, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 = 0. \quad (2)$$

On the other hand, if there exists $(\rho, M) \succeq B_W \in \mathbb{R}^d \times \mathbb{R}^d$ with $k M k_2 < 1$ such that the above equality holds, then ϕ must satisfy exact linear BC with $W = \frac{k \rho k_2}{1 - k M k_2}$.

Please see Appendix D for proof. The above shows that linear BC is equivalent to a simple linear relationship between the feature and the expected next step's feature and reward. This motivates our feature learning algorithm: if we are capable of learning a representation ϕ such that the transition

from the current feature $\phi(s, a)$ to the expected next feature $\mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e)]$ and reward $r(s, a)$ is linear, then we've found a feature ϕ that is linear BC.

4.1. Algorithm

To learn the representation that achieves linear BC, we use Proposition 4.2 to design a bilevel optimization program. We start with deterministic transition as a warm up.

Deterministic transition Due to determinism in the transition, we do not have an expectation with respect to s' anymore. So, we design the bilevel optimization as follows:

$$\begin{aligned} \min_{\phi \geq 0} \left[\min_{(\rho, M) \geq 0} \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \phi(s', \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right] \quad (3) \\ = \{(\rho, M) \geq B_{k\rho^*k} \mathbb{R}^{d \times d}; kMk_2 \leq kM^*k_2\}, \end{aligned}$$

where ρ^*, M^* are the optimal ρ, M for the linear BC ϕ^* (in Assumption 4.1). Namely, we aim to search for a representation $\phi \geq 0$, such that the relationship between $\phi(s, a)$ and the combination of the next time step's feature $\phi(s', \pi_e)$ and the reward, is linear. The spectral norm constraint in (3) is justified by Proposition 4.2. Solving the above bilevel moment condition finds a representation that achieves approximate linear BC.

However, there is no guarantee that such representation can provide a good coverage over π_e 's traces. We introduce regularizations for ϕ using the ideas from optimal designs, particularly the E -optimal design. Define the minimum eigenvalue regularization

$$R_E(\phi) := \lambda_{\min}(\mathbb{E}_D[\phi(s, a)\phi(s, a)^T]),$$

as the smallest eigenvalue of the empirical feature covariance matrix under the representation ϕ . Maximizing this quantity ensures that our feature provides good coverage, i.e. $\lambda_{\min}(\mathbb{E}_D[\phi(s, a)\phi(s, a)^T])$ is as large as possible.

Thus, to learn a representation that is approximately linear Bellman complete and also provides sufficient coverage, we formulate the following constrained bilevel optimization:

$$\begin{aligned} \min_{\phi \geq 0} \left[\min_{(\rho, M) \geq 0} \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \phi(s', \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right] \\ \text{s.t.}, R_E(\phi) \geq \beta/2. \end{aligned}$$

To extend this to stochastic transitions, there is an additional double sampling issue, which we discuss and address now.

Stochastic transition Ideally, we would solve the following bilevel optimization problem,

$$\min_{\phi \in \Phi} \left[\min_{(\rho, M) \in \Theta} \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e)] \\ r(s, a) \end{bmatrix} \right\|_2^2 \right]. \quad (4)$$

However, we cannot directly optimize the above objective since we do not have access to $P(s, a)$ to compute the expected next step feature $\mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e)]$. Also note that the expectation $\mathbb{E}_{s'}$ is inside the square which means that we cannot even get an unbiased estimate of the gradient of (ϕ, M) with one sample $s' \sim P(s, a)$. This phenomenon is related to the *double sampling* issue on offline policy evaluation literature which forbids one to directly optimize Bellman residuals. Algorithms such as TD are designed to overcome the double sampling issue. Here, we use a different technique to tackle this issue (Chen & Jiang, 2019). We introduce a function class $G \subseteq \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d$ which is rich enough to contain the expected next feature.

Assumption 4.3. For any $\phi \geq 0$, we have that the mapping $(s, a) \mapsto \gamma \mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e)]$ is in G .

We form the optimization problem as follows:

$$\begin{aligned} \min_{\phi \geq 0} \left[\min_{(\rho, M) \geq 0} \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \phi(s', \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right] \\ \min_{g \in G} \mathbb{E}_D \left[kg(s, a) - \gamma \phi(s', \pi_e) \right]^2. \quad (5) \end{aligned}$$

Note that under Assumption 4.3, the min over G will approximate $\gamma^2 \mathbb{E}_D k \mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e) - \phi(s', \pi_e)]^2$, i.e., the average variance induced by the stochastic transition. Thus, for a fixed ϕ and M , we can see that the following

$$\begin{aligned} \mathbb{E}_D \left[kM \phi(s, a) - \gamma \phi(s', \pi_e) \right]^2 \\ \gamma^2 \mathbb{E}_D k \mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e) - \phi(s', \pi_e)]^2 \end{aligned}$$

is indeed an unbiased estimate of:

$$\mathbb{E}_{s, a} \left\| M \phi(s, a) - \gamma \mathbb{E}_{s' \sim P(s, a)}[\phi(s', \pi_e)] \right\|_2^2$$

which matches to the ideal objective in Eq. 4. Thus solving for ϕ based Eq. 5 allows us to optimize Eq. 4, which allows us to learn an approximate linear Bellman complete representation. Similarly, we incorporate the E -optimal design here by adding a constraint that forcing the smallest eigenvalue of the empirical feature covariance matrix, i.e., $R_E(\phi)$, to be lower bounded.

Once we learn a representation that is approximately linear BC, and also induces sufficient coverage, we simply call the LSPE to estimate V^{π_e} . The whole procedure is summarized in Algorithm 2. Note in Alg 2 we put constraints to the objective function using Lagrangian multiplier.

There are other design choices that also encourage coverage. One particular choice we study empirically is motivated by the idea of D-optimal design. Here we aim to find a representation that maximizes the following log-determinant

$$R_D(\phi) := \ln \det(\mathbb{E}_D[\phi(s, a)\phi(s, a)^T]).$$

Algorithm 2 OPE with Bellman Complete and exploratory Representation Learning (BCRL)

- 1: **Input:** Representation class \mathcal{F} , dataset D of size $2N$, design regularization R , function class G , policy π_e .
- 2: Randomly split D into two sets D_1, D_2 of size N .
- 3: If the system is stochastic, learn representation $\hat{\phi}$ as,

$$\arg \min_{\phi \in \mathcal{F}} \left[\min_{(\rho, M) \in \mathcal{M}} \mathbb{E}_{D_1} \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma \phi(s^\theta, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right. \\ \left. \min_{g \in G} \mathbb{E}_{D_1} kg(s, a) - \gamma \phi(s^\theta, \pi_e) k_2^2 \right] + \lambda R(\phi)$$

- 4: Otherwise, for deterministic system, learn $\hat{\phi}$ as,

$$\arg \min_{\phi \in \mathcal{F}} \left[\min_{(\rho, M) \in \mathcal{M}} \mathbb{E}_{D_1} \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma \phi(s^\theta, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right. \\ \left. \lambda R(\phi) \right]$$

- 5: Return $\hat{V} := \text{LSPE}(\pi_e, \hat{\phi}, D_2)$.

When D is large, then the regularization $R_D(\phi)$ approximates $\sum_i \ln(\sigma_i(\mathbb{E}_\nu[\phi(s, a)\phi(s, a)^\top]))$. Maximizing $R_D(\phi)$ then intuitively maximizes coverage over all directions. D-optimal design is widely used for bandits (Lattimore & Szepesvári, 2020) and RL (Wang et al., 2021; Agarwal et al., 2019) to design exploration distributions with global coverage. Note that, in contrast to these contexts where the feature is fixed and the distribution is optimized, we optimize the feature, given the data distribution ν .

4.2. Sample Complexity Analysis

We now prove a finite sample utility guarantee for the empirical constrained bilevel optimization problem,

$$\hat{\phi} \geq \arg \min_{\phi \in \mathcal{F}} \left[\min_{(\rho, M) \in \mathcal{M}} \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma \phi(s^\theta, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \right. \\ \left. \min_{g \in G} \mathbb{E}_D kg(s, a) - \gamma \phi(s^\theta, \pi_e) k_2^2 \right]. \quad (6)$$

s.t., $R_E(\phi) \leq \beta/2$.

For simplicity, we state our results for discrete function class \mathcal{F} and G . Note that our sample complexity only scales with respect $\ln(j)$ and $\ln(jG)$, which are the standard complexity measures for discrete function classes. We extend our analysis to infinite function classes under metric entropy assumptions (Wainwright, 2019; van der Vaart & Wellner, 1996) in the Appendix; see Assumption E.5.

The following theorem shows that Algorithm 2 learns a representation $\hat{\phi}$ that is $O(N^{-1/2})$ approximately Linear BC and has coverage.

Theorem 4.4. Assume Assumption 4.1 (and Assumption 4.3 if the system is stochastic). Let $C_2 := \frac{96 \log^{1/2}(j) + 4\beta \bar{d} + 4 \log^{1/2}(1/\delta)}{\beta/4}$. If $N \geq C_2^2$, then for any $\delta \geq (0, 1)$, w.p. at least $1 - \delta$, we have

1. $\hat{\phi}$ satisfies $\hat{\epsilon}$ -approximate Linear BC with

$$\hat{\epsilon} \leq \frac{13d(1+W)^2 \log^{1/2}(4Wj) \log^{1/2}(jN/\delta)}{\beta N} + \frac{7\gamma(1+W) \log^{1/2}(2jGj/\delta)}{\beta N},$$

2. $\lambda_{\min}(\hat{\phi}) \geq \beta/4$.

If transitions are deterministic, treat $\log(jG) = 0$.

Proof Sketch. First, we use Weyl's Perturbation Theorem and chaining to show that the eigenvalues of $\hat{\phi}$ are close to $\hat{\phi}(\phi)$, uniformly over ϕ . This implies that (a) $\lambda_{\min}(\hat{\phi}(\phi^*)) \geq \beta/2$, and hence is feasible in the empirical bilevel optimization Equation (6), and (b) $\lambda_{\min}(\hat{\phi}(\phi)) \geq \beta/4$. Since ϕ^* is feasible, we apply uniform concentration arguments to argue that $\hat{\phi}$ has low population loss (Equation (5)), which implies approximate Linear BC. \square

The error term in $\hat{\epsilon}$ is comprised of the statistical errors of fitting \mathcal{M} and of fitting G for the double sampling correction. In the contextual bandit setting, i.e. $\gamma = 0$, there is no transition, so the second term becomes 0. Chaining together with Theorem 3.3 gives the following end-to-end $O(N^{-1/2})$ evaluation error guarantee for LSPE with the learned features:

Theorem 4.5. Under Assumption 4.1 (and Assumption 4.3 if $P(s, a)$ is stochastic). Let $C_2, \hat{\epsilon}$ be as defined in Theorem 4.4. If $N \geq C_2^2$, then we have for any $\delta \geq (0, 1/2)$, w.p. at least $1 - 2\delta$, for all distributions p_0 ,

$$\left| \mathbb{E}_{p_0}^{\pi_e} \mathbb{E}_s \left[\hat{f}_K(s, \pi_e) \right] \right| \leq \frac{\gamma^{K/2}}{1 - \gamma} + \frac{4 \sqrt{\left\| \frac{dd_{p_0}^{\pi_e}}{d\nu} \right\|_1}}{(1 - \gamma)^2} \hat{\epsilon} \\ + \frac{960 \beta^{-1/2} (1+W) d \log(N) \sqrt{\log(10/\delta)}}{(1 - \gamma)^2 \beta N}.$$

Comparison to FQE What if one ignores the representation learning and just runs the Fitted Q-Evaluation (FQE) which directly performs least square fitting with the nonlinear function class $F := \{f_w \succ \phi(s, a) : \phi \in \mathcal{F}, kwk_2 \leq W\}$? As $N \rightarrow \infty$, FQE will suffer the following worst-case Bellman error (also called inherent Bellman error):

$$\epsilon_{\text{ibe}} := \max_{f \in F} \min_{g \in G} kg - \mathbb{E}^{\pi_e} f k_\nu^2.$$

Algorithm 3 Practical Instantiation of BCRL

- 1: **Input:** Offline dataset $D = \bar{f}_s, a, s^{\bar{\rho}}, g$, target policy π_e
- 2: Initialize parameters for ϕ , M , ρ , and $\bar{\phi}$
- 3: **for** $t = 0, 1, \dots, T$ **do**
- 4: $M_{t+1} \leftarrow M_t - \eta \Gamma_M J(\phi_t, M_t, \rho_t, \bar{\phi}_t)$
- 5: $\rho_{t+1} \leftarrow \rho_t - \eta \Gamma_\rho J(\phi_t, M_t, \rho_t, \bar{\phi}_t)$
- 6: $\phi_{t+1} \leftarrow \phi_t - \eta \Gamma_\phi J(\phi_t, M_{t+1}, \rho_{t+1}, \bar{\phi}_t)$
- 7: $\bar{\phi}_{t+1} \leftarrow \tau \phi_{t+1} + (1 - \tau) \bar{\phi}_t$
- 8: **end for**
- 9: Linear evaluation: $\hat{V} = \text{LSPE}(\pi_e, \phi_T, D)$.

Note that our assumption that there exists a linear BC representation ϕ^* does not imply that the worst-case Bellman error ε_{ibe} is small. In contrast, when $N \rightarrow \infty$, our approach will accurately estimate $V_{\rho_0}^{\pi_e}$.

5. A Practical Implementation

In this section we instantiate a practical implementation to learn our representation using deep neural networks for our representation function class ϕ . Based on Equation (3), we first formalize our bilevel optimization objective:

$$J(\phi, M, \rho, \bar{\phi}) = \mathbb{E}_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \bar{\phi}(s^{\bar{\rho}}, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \\ \lambda \log \det \mathbb{E}_D [\phi(s, a) \phi(s, a)^\top].$$

In our implementation, we replace the hard constraint presented in Section 4.1 with a Lagrange multiplier, i.e. we use the optimal design constraint as a regularization term when learning ϕ . Specifically, we maximize the log det of the covariance matrix induced by the feature, which maximizes all eigenvalues since log det is the sum of the log eigenvalues. Our experiment results in Section 6.4 demonstrate that it indeed improves the condition number of ϕ . In some of our experiments, we use a target network in our implementation. Namely, the updates make use of a target network $\bar{\phi}$ where the weights can be an exponential moving average of the representation network’s weights. This idea of using a slow moving target network has been shown to stabilize training in both the RL (Mnih et al., 2013) and the representation learning literature (Grill et al., 2020).

As summarized in Algorithm 3, given our offline behavior dataset D and target policy π_e , we iteratively update M

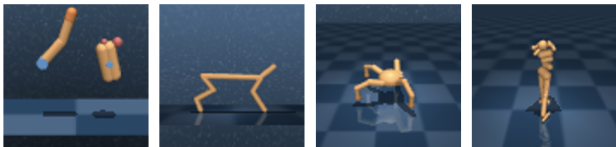


Figure 1. Representative frames from DeepMind Control Suite tasks. From left to right, Finger Turn Hard, Cheetah Run, Quadruped Walk, and Humanoid Stand.

and ϕ and then use the resulting learned representation to perform OPE. Please see Appendix F for implementation and hyperparameter details. As we will show in our experiments, our update procedures for ϕ significantly minimizes the Bellman Completeness loss and also improve the condition number of ϕ , which are the two key quantities to ensure good performance of LSPE with linear regression as shown in Theorem 3.3.

6. Experiments

Our goal is to answer the following questions. (1) How do our representations perform on downstream LSPE compared to other popular unsupervised representation learning techniques? (2) How important are both the linear bellman completeness and optimal design components for learning representations? (3) How competitive is BCRL with other OPE methods, especially for evaluating beyond the initial state distribution?

Following the standards in evaluating representations in supervised learning, we used a **linear evaluation protocol**, i.e., linear regression in LSPE on top of a given representation. This allows us to focus on evaluating the quality of the representation. We compared our method to prior techniques on a range of challenging, image-based continuous control tasks from the DeepMind Control Suite benchmark (Tassa et al., 2018): Finger Turn Hard, Cheetah Run, Quadruped Walk, and Humanoid Stand. Frames from the tasks are shown in Figure 1. To investigate our learned representation, we benchmark our representation against two state-of-the-art self-supervised representation learning objectives adopted for RL: (1) CURL uses the InfoNCE objective to contrastively learn state-representations; and (2) SPR adopts the BYOL contrastive learning framework to learn representations with latent dynamics. Note, we modified CURL for OPE by optimizing the contrastive loss between state-action pairs rather than just states. For SPR, we did not include the Q prediction head and used their state representation plus latent dynamics as the state-action representation for downstream linear evaluation. We used the same architecture for the respective representations across all evaluated algorithms. To compare against other OPE methods, we additionally compared against Fitted Q-Evaluation (Munos & Szepesvári, 2008; Kostrikov & Nachum, 2020) (FQE), weighted doubly robust policy evaluation (Jiang & Li, 2016; Thomas & Brunskill, 2016) (DR), Dreamer-v2 (Hafner et al., 2021) model based evaluation (MB), and DICE (Yang et al., 2020). We modified implementations from the benchmark library released by Fu et al. (2021) for FQE and DR and used the authors’ released codebases for Dreamer-v2 and BestDICE.²

²https://github.com/google-research/dice_rl.

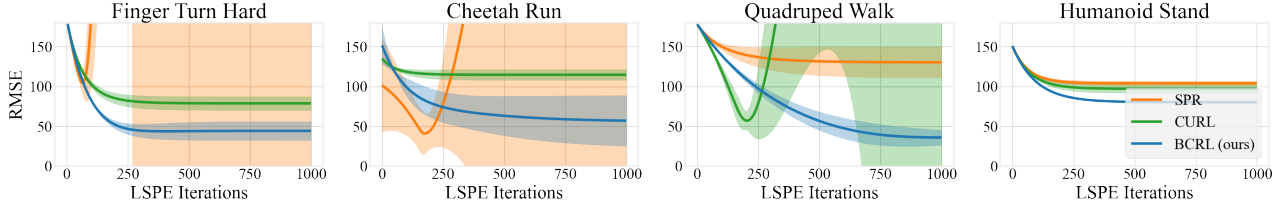


Figure 2. OPE curves across five seeds using representations trained with BCRL, SPR, and CURL on the offline datasets (Table 1).

Our target policies were trained using the authors’ implementation of DRQ-v2 (Yarats et al., 2021), a state-of-the-art, off-policy actor critic algorithm for vision-based continuous control. With high-quality target policies, we collected 200 rollouts and did a linear evaluation protocol to predict the discounted return. For our offline behavior datasets, we collected 10K samples from a trained policy with mean performance roughly a quarter of that of the target policy (Table 1). All results are aggregated over five random seeds. See Appendix F for details on hyperparameters, environments, and dataset composition.

6.1. OPE via LSPE with Learned Representations

Figure 2 compares the OPE performance of BCRL against SPR and CURL. Representations learned by BCRL outperform those learned by SPR and CURL. On some tasks, SPR and CURL both exhibited an exponential error amplification with respect to the number of iterations of LSPE, while BCRL did not suffer from any blowup.

6.2. OPE Performance

Figure 3 compares the OPE performance of BCRL against multiple benchmarks from the OPE literature. BCRL is competitive with FQE and evaluates better than other benchmarks across the tasks that we tested on.

We also evaluated how well estimated values from BCRL rank policies. Following (Fu et al., 2021), we use the spearman ranking correlation metric to compute the correlation between ordinal rankings according to the OPE estimates and the ground truth values. For ranking, we evaluated three additional target policies with mean performances roughly 75%, 50%, and 10% of the target policy. Figure 3 presents the mean correlation of each evaluation algorithm across all tasks. BCRL is competitive with FQE and consistently better than other benchmarks at ranking policies.

6.3. Further Investigation of Different Settings

In this section, we consider two additional settings: (1) the offline dataset contains some on-policy data, which ensures that the offline data provides coverage over the evaluation policy’s state action distribution; (2) we evaluate all esti-

mators beyond the original initial state distribution d_0 . The first setting ensures that baselines and our algorithm all satisfy the coverage condition, thus demonstrating the unique benefit of learning a Bellman complete representation. The second setting evaluates the robustness of our algorithm, i.e., the ability to estimate beyond the original d_0 .

Evaluation on On-Policy + Off-Policy Data: To further investigate the importance of learning linear BC features, we experiment with learning representations from an offline dataset that also contains state-action pairs from the target policy. More specifically, we train our representations on a dataset containing 10K behavior policy and 10K target policy samples. Note, only the experiment in this paragraph uses this mixture dataset. With the addition of on-policy data from the target policy, we can focus on just the role of linear BC for OPE performance because the density ratio $\frac{d_{\pi_e}}{d_{\pi}}$ and the relative condition number (Eq. 1) is at most 2, i.e. we omit the design regularization and focused on minimizing the Bellman completeness loss. Figure 4 (Left) shows that BCRL outperforms baselines in this setting, even

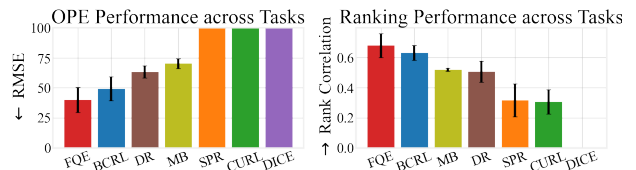


Figure 3. (Left) Root mean squared evaluation error across all tasks. (Right) Mean spearman ranking correlation across all tasks.

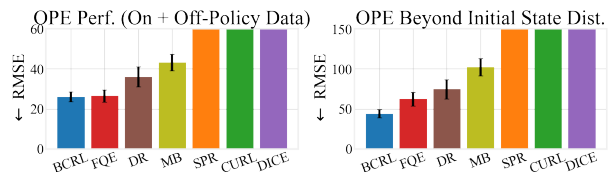


Figure 4. (Left) Evaluation on a mixture dataset with on-policy and off-policy data. Note the addition of target policy data bounds the relative condition number (Eq. 1). (Right) Evaluation beyond the initial state distribution (given just the offline data).

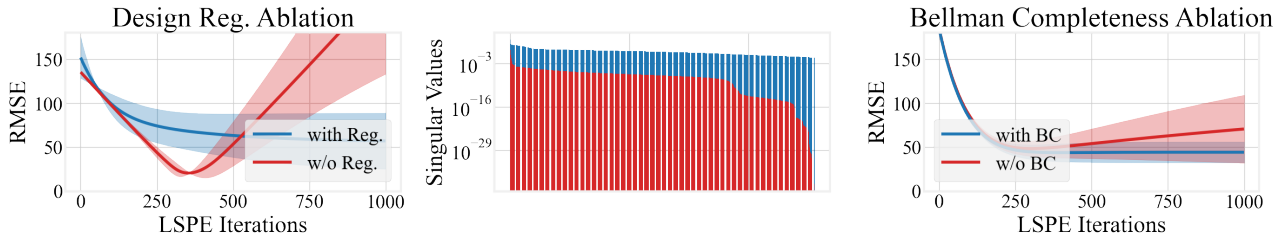


Figure 5. (Left) BCRL’s OPE curves for Cheetah Run with (blue) and without (red) the D -optimal design based regularization. (Center) Bar plot of the singular values of the feature covariance matrices for the left plot. (Right) OPE curves for Finger Turn Hard with (blue) and without (red) optimizing for linear BC.

matching FQE performance across tasks. Our experiments corroborate our analysis that explicitly enforcing our learned representations to be linear BC improves OPE. Note that the fact that LSPE with SPR and CURL features still blows up under this mixture data means that the other representations’ failures are not just due to coverage.

Evaluation Beyond Initial State Distribution: To further investigate the benefits of optimizing the D -optimal design to improve coverage, we investigate doing OPE *beyond* the initial state distribution d_0 , which is supported by Theorem 4.5. Note that if our representation is exactly Bellman complete and also the corresponding feature covariance matrix is well-conditioned, we should be able to evaluate well on any states. Specifically, we experiment on evaluating at all timesteps in a target policy rollout, not just at the initial state distribution. Figure 4 (Right) shows that BCRL is able to more robustly evaluate out-of-distribution than all other benchmarks.

6.4. Ablation Studies

Impact of Optimal Design Regularization: To investigate the impact of maximizing the D -optimal design, we ablate the design regularization term from our objective and analyze the downstream evaluation performance and the respective feature covariance matrices on the offline dataset. Figure 5 (Center) presents a bar plot of the singular values of the feature covariance matrix ($\Sigma(\phi) := \mathbb{E}_{s,a} \nu \phi(s,a)\phi(s,a)^\top$). Figure 5 (Left) shows the downstream OPE performance for features trained with and without the design regularization on the Cheetah Run task. Note that without the regularization, we find that the feature covariance matrix has much worse condition number, i.e. feature is less exploratory. As our analysis suggests, we also observe a deterioration in evaluation performance without the design regularization to explicitly learn exploratory features.

Impact of Linear Bellman Completeness: Figure 5 (Right) presents an ablation study where we only optimize for the design term in our objective. We find that

downstream OPE performance degrades without directly optimizing for linear BC, suggesting that a feature with good coverage alone is not enough to avoid error amplification.

7. Conclusion

We present BCRL which leverages rich function approximation to learn Bellman complete and exploratory representations for stable and accurate offline policy evaluation. We provide a mathematical framework of representation learning in offline RL, which generalizes all existing representation learning frameworks from the RL theory literature. We provide an end-to-end theoretical analysis of our approach for OPE and demonstrate that BCRL can accurately estimate policy values with polynomial sample complexity. Notably, the complexity has no explicit dependence on the size of the state and action space, instead, it only depends on the statistical complexity of the representation hypothesis class. Experimentally, we extensively evaluate our approach on the DeepMind Control Suite, a set of image-based, continuous control, robotic tasks. First, we show that under the linear evaluation protocol – using linear regression on top of the representations inside the classic LSPE framework – our approach outperforms prior RL representation techniques CURL and SPR which leverage contrastive representation learning techniques SimCLR and BYOL respectively. We also show that BCRL achieves competitive OPE performance with the state-of-the-art FQE, and noticeably improves upon it when evaluating beyond the initial state distribution. Finally, our ablations show that approximate Linear Bellman Completeness and coverage are crucial ingredients to the success of our algorithm. Future work includes extending BCRL to offline policy optimization.

ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation under Grant No. 1846210 and by a Cornell University Fellowship. We thank Rahul Kidambi, Ban Kawas, and the anonymous reviewers for useful discussions and feedback.

References

- Agarwal, A., Jiang, N., Kakade, S. M., and Sun, W. Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 2019.
- Agarwal, A., Kakade, S., Krishnamurthy, A., and Sun, W. Flambe: Structural complexity and representation learning of low rank mdps. *NeurIPS*, 33:20095–20107, 2020.
- Amortila, P., Jiang, N., and Xie, T. A variant of the wang-foster-kakade lower bound for the discounted setting. *arXiv preprint arXiv:2011.01075*, 2020.
- Anand, A., Racah, E., Ozair, S., Bengio, Y., Côté, M., and Hjelm, R. D. Unsupervised state representation learning in atari. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), *NeurIPS*, pp. 8766–8779, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/6fb52e71b837628ac16539c1ff911667-Abstract.html>.
- Antos, A., Szepesvári, C., and Munos, R. Fitted q-iteration in continuous action-space mdps. In Platt, J., Koller, D., Singer, Y., and Roweis, S. (eds.), *NIPS*, volume 20. Curran Associates, Inc., 2008. URL <https://proceedings.neurips.cc/paper/2007/file/da0d1111d2dc5d489242e60ebcbaf988-Paper.pdf>.
- Bhatia, R. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
- Boucheron, S., Lugosi, G., and Massart, P. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Chen, J. and Jiang, N. Information-theoretic considerations in batch reinforcement learning. In *ICML*, pp. 1042–1051. PMLR, 2019.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. E. A simple framework for contrastive learning of visual representations. *CoRR*, abs/2002.05709, 2020. URL <https://arxiv.org/abs/2002.05709>.
- Chung, W., Nath, S., Joseph, A., and White, M. Two-timescale networks for nonlinear value function approximation. In *ICLR*, 2019. URL <https://openreview.net/forum?id=rJleN20qK7>.
- Du, S. S., Kakade, S. M., Lee, J. D., Lovett, S., Mahajan, G., Sun, W., and Wang, R. Bilinear classes: A structural framework for provable generalization in rl. *arXiv preprint arXiv:2103.10897*, 2021.
- Duan, Y., Jia, Z., and Wang, M. Minimax-optimal off-policy evaluation with linear function approximation. In *ICML*, pp. 2701–2709. PMLR, 2020.
- Foster, D. J., Krishnamurthy, A., Simchi-Levi, D., and Xu, Y. Offline reinforcement learning: Fundamental barriers for value function approximation. *arXiv preprint arXiv:2111.10919*, 2021.
- Fu, J., Norouzi, M., Nachum, O., Tucker, G., Wang, Z., Novikov, A., Yang, M., Zhang, M. R., Chen, Y., Kumar, A., Paduraru, C., Levine, S., and Paine, T. L. Benchmarks for deep off-policy evaluation. *CoRR*, abs/2103.16596, 2021. URL <https://arxiv.org/abs/2103.16596>.
- Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., Piot, B., kavukcuoglu, k., Munos, R., and Valko, M. Bootstrap your own latent - a new approach to self-supervised learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *NeurIPS*, volume 33, pp. 21271–21284. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf>.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, pp. 1861–1870. PMLR, 2018.
- Hafner, D., Lillicrap, T. P., Norouzi, M., and Ba, J. Mastering atari with discrete world models. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=0oabwyZbOu>.
- Hao, B., Duan, Y., Lattimore, T., Szepesvári, C., and Wang, M. Sparse feature selection makes batch reinforcement learning more sample efficient. In *ICML*, pp. 4063–4073. PMLR, 2021.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. B. Momentum contrast for unsupervised visual representation learning. *CoRR*, abs/1911.05722, 2019. URL <http://arxiv.org/abs/1911.05722>.
- Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. In *ICML*, pp. 652–661. PMLR, 2016.
- Jiang, N., Krishnamurthy, A., Agarwal, A., Langford, J., and Schapire, R. E. Contextual decision processes with low bellman rank are pac-learnable. In *ICML*, pp. 1704–1713. PMLR, 2017.

- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *COLT*, pp. 2137–2143. PMLR, 2020.
- Kostrikov, I. and Nachum, O. Statistical bootstrapping for uncertainty estimation in off-policy evaluation, 2020. URL <https://arxiv.org/abs/2007.13609>.
- Laskin, M., Srinivas, A., and Abbeel, P. Curl: Contrastive unsupervised representations for reinforcement learning. In *ICML*, pp. 5639–5650. PMLR, 2020.
- Lattimore, T. and Szepesvári, C. *Bandit algorithms*. Cambridge University Press, 2020.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In *ICLR*, 2016.
- Mazouze, B., Tachet des Combes, R., Doan, T. L., Bachman, P., and Hjelm, R. D. Deep reinforcement and infomax learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *NeurIPS*, volume 33, pp. 3686–3698. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/26588e932c7ccfa1df309280702fe1b5-Paper.pdf>.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. A. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>.
- Modi, A., Chen, J., Krishnamurthy, A., Jiang, N., and Agarwal, A. Model-free representation learning and exploration in low-rank mdps. *arXiv preprint arXiv:2102.07035*, 2021.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *JMLR*, 9(5), 2008.
- Nachum, O. and Yang, M. Provable representation learning for imitation with contrastive fourier features. *CoRR*, abs/2105.12272, 2021. URL <https://arxiv.org/abs/2105.12272>.
- Nedić, A. and Bertsekas, D. P. Least squares policy evaluation algorithms with linear function approximation. *Discrete Event Dynamic Systems*, 13(1):79–110, 2003.
- Ni, C., Zhang, A. R., Duan, Y., and Wang, M. Learning good state and action representations via tensor decomposition. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pp. 1682–1687. IEEE, 2021.
- Parr, R., Li, L., Taylor, G., Painter-Wakefield, C., and Littman, M. L. An analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *ICML*, pp. 752–759, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi: 10.1145/1390156.1390251. URL <https://doi.org/10.1145/1390156.1390251>.
- Pollard, D. Empirical processes: theory and applications. In *NSF-CBMS regional conference series in probability and statistics*, pp. i–86. JSTOR, 1990.
- Schwarzer, M., Anand, A., Goel, R., Hjelm, R. D., Courville, A. C., and Bachman, P. Data-efficient reinforcement learning with momentum predictive representations. *CoRR*, abs/2007.05929, 2020. URL <https://arxiv.org/abs/2007.05929>.
- Song, Z., Parr, R. E., Liao, X., and Carin, L. Linear feature encoding for reinforcement learning. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R. (eds.), *NIPS*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/8232e119d8f59aa83050a741631803a6-Paper.pdf>.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- Thomas, P. S. and Brunskill, E. Data-efficient off-policy policy evaluation for reinforcement learning, 2016. URL <https://arxiv.org/abs/1604.00923>.
- Uehara, M. and Sun, W. Pessimistic model-based offline reinforcement learning under partial coverage. *arXiv preprint arXiv:2107.06226*, 2021.
- Uehara, M., Zhang, X., and Sun, W. Representation learning for online and offline RL in low-rank MDPs. In *ICLR*, 2022. URL <https://openreview.net/forum?id=J4iSIR9fhY0>.
- van den Oord, A., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. URL <http://arxiv.org/abs/1807.03748>.
- van der Vaart, A. W. and Wellner, J. A. *Weak Convergence and Empirical Processes*. Springer Series in Statistics. Springer New York, 1996. ISBN 9781475725476. doi: 10.1007/978-1-4757-2545-2. URL <http://link.springer.com/10.1007/978-1-4757-2545-2>.

- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Wang, R., Foster, D. P., and Kakade, S. M. What are the statistical limits of offline RL with linear function approximation? *CoRR*, abs/2010.11895, 2020. URL <https://arxiv.org/abs/2010.11895>.
- Wang, Y., Wang, R., and Kakade, S. M. An exponential lower bound for linearly-realizable mdps with constant suboptimality gap. *arXiv preprint arXiv:2103.12690*, 2021.
- Yang, M. and Nachum, O. Representation matters: Offline pretraining for sequential decision making. In Meila, M. and Zhang, T. (eds.), *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11784–11794. PMLR, 2021. URL <http://proceedings.mlr.press/v139/yang21h.html>.
- Yang, M., Nachum, O., Dai, B., Li, L., and Schuurmans, D. Off-policy evaluation via the regularized lagrangian. *NeurIPS*, 33:6551–6561, 2020.
- Yarats, D., Zhang, A., Kostrikov, I., Amos, B., Pineau, J., and Fergus, R. Improving sample efficiency in model-free reinforcement learning from images. *CoRR*, abs/1910.01741, 2019. URL <http://arxiv.org/abs/1910.01741>.
- Yarats, D., Fergus, R., Lazaric, A., and Pinto, L. Mastering visual continuous control: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021.
- Zhang, W., He, J., Zhou, D., Zhang, A., and Gu, Q. Provably efficient representation learning in low-rank markov decision processes. *arXiv preprint arXiv:2106.11935*, 2021.

Appendices

A. Metric Entropy and Entropy Integral

We recall the standard notions of entropy integrals here, based on the following distance on \mathcal{F} ,

$$d(\phi, \tilde{\phi}) = \mathbb{E}_D \left[\left\| \phi(s, a) - \tilde{\phi}(s, a) \right\|_2 \right]$$

Let $N(t, \mathcal{F})$ denote the t -covering number under d .

Definition A.1. Define the entropy integral, which we assume to be finite as,

$$\kappa(\mathcal{F}) = \int_0^4 \log^{1/2} N(t, \mathcal{F}) dt$$

When \mathcal{F} is finite, $N(t) \leq |\mathcal{F}|$, so $\kappa(\mathcal{F}) \leq O(\log^{1/2} |\mathcal{F}|)$.

B. Technical Lemmas

Lemma B.1. Let X_i be i.i.d. random variables s.t. $|X_i| \leq c$ and $\mathbb{E}[X_i^2] \leq \nu$, then for any $\delta \in (0, 1)$, we have w.p. at least $1 - \delta$,

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X] \right| \leq \inf_{a>0} \frac{\nu}{2a} + \frac{(c+a) \log(2/\delta)}{N},$$

and if $\nu \leq L\mathbb{E}[X]$ for some positive L , then in particular,

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X] \right| \leq \frac{1}{2} \mathbb{E}[X] + \frac{(c+L) \log(2/\delta)}{N}.$$

Proof. First, by Bernstein's inequality (Boucheron et al., 2013, Theorem 2.10), we have w.p. $1 - \delta$,

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mathbb{E}[X] \right| \leq \sqrt{\frac{2\nu \log(2/\delta)}{N}} + \frac{c \log(2/\delta)}{N}$$

Using the fact that $2xy \leq x^2/a + ay^2$ for any $a > 0$, split the square root term,

$$\inf_{a>0} \frac{\nu}{2a} + \frac{(c+a) \log(2/\delta)}{N}$$

which yields the first part. If $\nu \leq L\mathbb{E}[X]$, picking $a = L$ concludes the proof. \square

We now state several results of Orlicz norms (mostly from Pollard, 1990) for completeness. For an increasing, convex, positive function $\psi: \mathbb{R}^+ \rightarrow \mathbb{R}^+$, such that $\psi(x) \geq [0, 1]$, define the Orlicz norm as

$$\|Z\|_{\psi} := \inf \{C > 0 \mid \mathbb{E}[\psi(|Z|/C)] \leq 1\}.$$

It is indeed a norm on the ψ -Orlicz space of random variables $L_{\psi}(\nu)$, since it can be interpreted as the Minkowski functional of the convex set $K = \{fX : \mathbb{E}[\psi(|fX|)] \leq 1\}$.

Let x_1, x_2, \dots, x_N be i.i.d. datapoints drawn from some underlying distribution. Let ω denote the randomness of the N sampled datapoints, and let $F_{\omega} = \{(f(x_i(\omega)))_{i=1}^N : f \in \mathcal{F}\} \subset \mathbb{R}^N$ denote the (random) set of vectors from the data corresponding to ω . Let σ denote a vector of N i.i.d. Rademacher random variables (± 1 equi-probably), independent of all else.

Lemma B.2 (Symmetrization). *For any increasing, convex, positive function ψ , we have*

$$\mathbb{E}_\omega \left[\left(\sup_{f \in \mathcal{F}_\omega} \left| \sum_{i=1}^N f_i - \mathbb{E}_\omega f_i \right| \right) \right] \leq \mathbb{E}_\omega \left[\left(2 \sup_{f \in \mathcal{F}_\omega} \psi(f) \right) \right].$$

Proof. See Theorem 2.2 of (Pollard, 1990). □

Lemma B.3 (Contraction). *Let $\mathcal{F} \subseteq \mathbb{R}^N$, and suppose $\lambda : \mathbb{R}^N \rightarrow \mathbb{R}^N$ s.t. each component $\lambda_i : \mathbb{R} \rightarrow \mathbb{R}$ is L -Lipschitz. Then, for any increasing, convex, positive function ψ , we have*

$$\mathbb{E} \left[\left(\sup_{f \in \mathcal{F}} \psi(\lambda(f)) \right) \right] \leq \frac{3}{2} \mathbb{E} \left[\left(2L \sup_{f \in \mathcal{F}} \psi(f) \right) \right]$$

Proof. Apply Theorem 5.7 of (Pollard, 1990) to the functions λ_i/L , which are contractions. □

We now focus on the specific Orlicz space of sub-Gaussian random variables, with the function $\psi(x) = \frac{1}{5} \exp(x^2)$. The ψ -Orlicz norm of the maximum of random variables can be bounded by the maximum of the ψ -Orlicz norms.

Lemma B.4. *For any random variables Z_1, \dots, Z_m , we have*

$$\max_{i=1, \dots, m} \|Z_i\|_{\psi} \leq \sqrt{2 + \log(m)} \max_{i=1, \dots, m} \|Z_i\|_{\psi}$$

Proof. See Lemma 3.2 of (Pollard, 1990). □

Lemma B.5. *For each $f \in \mathbb{R}^N$, we have $\|f\|_{\psi} \leq 2\|f\|_2$.*

Proof. See Lemma 3.1 of (Pollard, 1990). □

The following is a truncated chaining result for Orlicz norms. This result is not new, but many sources state and prove it in terms of covering and Rademacher complexity, rather than for Orlicz norms. In particular, it generalizes Theorem 3.5 of (Pollard, 1990) – consider a sequence of α 's converging to zero.

Lemma B.6 (Chaining). *Let $\mathcal{F} \subseteq \mathbb{R}^N$ such that $0 \in \mathcal{F}$. Then,*

$$\left\| \sup_{f \in \mathcal{F}} \psi(f) \right\| \leq \inf_{\alpha > 0} \left\{ 3\alpha \sqrt{N} + 9 \int_{\alpha}^b \sqrt{\log(D(\delta/2, \mathcal{F}))} d\delta \right\},$$

where $b = \sup_{f \in \mathcal{F}} \|f\|_2$ and $D(\delta, \mathcal{F})$ is the Euclidean δ -packing number for \mathcal{F} .

Proof. Suppose b and all the packing numbers are finite, otherwise the right hand side is infinite and there is nothing to show. For an arbitrary $K > 1$, construct a sequence of K finer and finer approximations to \mathcal{F} ,

$$\mathcal{F}_0 = \mathcal{F} \quad \mathcal{F}_1 \quad \dots \quad \mathcal{F}_K \quad \mathcal{F}_{K+1} = \mathcal{F}$$

where for any $k \in [K]$, \mathcal{F}_k satisfies the property that for any $f \in \mathcal{F}$, there exist $n_k(f) \in \mathcal{F}_k$ s.t. $\|n_k(f) - f\|_2 \leq b 2^{-k}$. Indeed this can be done iteratively: for any \mathcal{F}_k , we can construct \mathcal{F}_{k+1} by adding elements to construct a maximal $b 2^{-k}$ packing (maximality ensures the distance requirement, since the existence of any vector which has larger distance can be added to the packing). By definition of $D(\cdot, \mathcal{F})$, we have $|\mathcal{F}_{k+1}| \leq D(b 2^{-k}, \mathcal{F})$.

For any $k \geq [K]$, we have by triangle inequality,

$$\begin{aligned} & \sup_{f \in F_{k+1}} j\sigma(f) - \sup_{f \in F_k} j\sigma(f) + \sup_{f \in F_{k+1}} j\sigma(n_k(f)) - \sup_{f \in F_k} j\sigma(n_k(f)) \\ &= \sup_{f \in F_k} j\sigma(f) - \sup_{f \in F_{k+1}} j\sigma(f) + \sup_{f \in F_{k+1}} j\sigma(n_k(f)) - \sup_{f \in F_{k+1}} j\sigma(f) \end{aligned}$$

If $k = K$, we can loosely bound the right-most term by Cauchy-Schwarz, since for any $f \in F$, we have $j\sigma(n_k(f)) - j\sigma(f) \leq \frac{1}{\sqrt{N}} \|k\sigma(n_k(f)) - k\sigma(f)\|_2 \leq \frac{1}{\sqrt{N}} \sqrt{2} \|k\sigma(n_k(f)) - k\sigma(f)\|_2$. So, we have

$$\left\| \sup_{f \in F} j\sigma(f) \right\| \leq \left\| \max_{f \in F_K} j\sigma(f) \right\| + \frac{\rho \sqrt{N} b^2}{\log 5} K,$$

since for any non-negative constant c , $\frac{1}{\sqrt{N}} \sqrt{2} \|k\sigma(n_k(f)) - k\sigma(f)\|_2 = \inf \{C > 0 : \frac{1}{5} \exp((c/C)^2) - 1\} = \frac{\rho c}{\log 5}$. If $k < K$, the suprema are taken over finite sets, so the maximum is attained. Hence, we can apply a special property of the ρ -Orlicz norm (Lemma B.4), to get,

$$\begin{aligned} \left\| \max_{f \in F_{k+1}} j\sigma(f) \right\| &\leq \left\| \max_{f \in F_k} j\sigma(f) \right\| + \left\| \max_{f \in F_{k+1}} j\sigma(n_k(f)) - j\sigma(f) \right\| \\ &\leq \left\| \max_{f \in F_k} j\sigma(f) \right\| + \sqrt{2 + \log(jF_{k+1})} \max_{f \in F_{k+1}} \|k\sigma(n_k(f)) - k\sigma(f)\|_2 \end{aligned}$$

By Lemma B.5,

$$\begin{aligned} & \left\| \max_{f \in F_k} j\sigma(f) \right\| + 2\sqrt{2 + \log(jF_{k+1})} \max_{f \in F_{k+1}} \|k\sigma(n_k(f)) - k\sigma(f)\|_2 \\ & \leq \left\| \max_{f \in F_k} j\sigma(f) \right\| + 2\sqrt{2 + \log(jF_{k+1})} b^2 k. \end{aligned}$$

Also, note that since $F_0 = \{0\}$ by construction, we have $\max_{f \in F_0} j\sigma(f) = 0$. Unrolling this, we have

$$\max_{f \in F} j\sigma(f) \leq \frac{\rho \sqrt{N} b^2}{\log 5} K + \sum_{k=0}^{K-1} 2^k b^2 \sqrt{2 + \log(D(b^2 2^{k+1}), F)}$$

Since for any $D \geq 2$, we have $\sqrt{2 + \log(1 + D)} \leq \frac{\rho}{D}$, 2.2,

$$\begin{aligned} & \frac{\rho \sqrt{N} b^2}{\log 5} K + 4.4b \sum_{k=0}^{K-1} 2^k \sqrt{\log(D(b^2 2^{k+1}), F)} \\ &= \frac{\rho \sqrt{N} b^2}{\log 5} K + 17.6b \sum_{k=0}^{K-1} (2^{k+1})^{\frac{1}{2}} (2^{k+2})^{\frac{1}{2}} \sqrt{\log(D(b^2 2^{k+1}), F)} \end{aligned}$$

Since $D(\cdot, F)$ is a monotone decreasing,

$$\begin{aligned} & \frac{\rho \sqrt{N} b^2}{\log 5} K + 17.6 \int_{b^2 2^{K+1}}^{b^2} \sqrt{\log(D(\delta, F))} d\delta \\ &= \frac{\rho \sqrt{N} b^2}{\log 5} K + 8.8 \int_{b^2 2^K}^b \sqrt{\log(D(\delta/2, F))} d\delta. \end{aligned}$$

Now consider any $\alpha > 0$. Pick K such that $\frac{b}{2^{K+1}} \leq \alpha \leq \frac{b}{2^K}$, then we have

$$\max_{f \in F} j\sigma(f) \leq \frac{2\alpha \rho \sqrt{N}}{\log 5} + 8.8 \int_{\alpha}^b \sqrt{\log(D(\delta/2, F))} d\delta.$$

Since α was arbitrary, the above bound holds when we take an infimum over α . □

C. Proofs for LSPE

We first show a generalization of the performance difference lemma (PDL).

Lemma C.1 (Generalized PDL). *For any policies π, π^θ , any function $f : S \times A \rightarrow \mathbb{R}$, and any initial state distribution μ , we have*

$$V_\mu^\pi - E_{s \sim \mu} [f(s, \pi^\theta)] = \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [T^{\pi^\theta} f(s, a) - f(s, \pi^\theta)] \quad (7)$$

Proof. Let T^π be the distribution of trajectories $\tau = (s_0, a_0, s_1, a_1, s_2, a_2, \dots)$ from rolling out π . So, we have,

$$\begin{aligned} V_\mu^\pi - E_{s \sim \mu} [f(s, \pi^\theta)] &= E_\tau [T^\pi \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] - E_{s \sim \mu} [f(s, \pi^\theta)]] \\ &= E_\tau [T^\pi \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + f(s_t, \pi^\theta) - f(s_t, \pi^\theta)) - f(s_0, \pi^\theta) \right]] \\ &= E_\tau [T^\pi \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \gamma f(s_{t+1}, \pi^\theta) - f(s_t, \pi^\theta)) \right]] \\ &= \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [r(s, a) + \gamma E_{s^\theta \sim P(s, a)} [f(s^\theta, \pi^\theta)] - f(s, \pi^\theta)] \\ &= \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [T^{\pi^\theta} f(s, a) - f(s, \pi^\theta)]. \end{aligned}$$

□

This generalizes the PDL, which we can get by setting $f(s, a) = Q^{\pi^\theta}(s, a)$:

$$\begin{aligned} V_\mu^\pi - V_\mu^{\pi^\theta} &= \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [T^{\pi^\theta} Q^{\pi^\theta}(s, a) - Q^{\pi^\theta}(s, \pi^\theta)] \\ &= \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [Q^{\pi^\theta}(s, a) - V^{\pi^\theta}(s)] \\ &= \frac{1}{1-\gamma} E_{s, a \sim d_\mu^\pi} [A^{\pi^\theta}(s, a)]. \end{aligned}$$

To prove our LSPE guarantee, we'll instantiate $f(s, a)$ to be the estimated $\hat{f}(s, a)$ from LSPE, and set $\pi = \pi^\theta$. This gives us an expression for the prediction error of LSPE,

$$\left| V_\mu^\pi - E_{s \sim \mu} [\hat{f}_k(s, \pi)] \right| = \frac{1}{1-\gamma} \left| E_{d_\mu^\pi} [T^\pi \hat{f}_k(s, a) - f_k(s, a)] \right|.$$

We then upper bound the right hand side by its $L_2(d_\mu^\pi)$ norm, which is the Bellman error. Next, we bound the Bellman error of running LSPE by the regression losses at each step.

Lemma C.2. *Consider any policy π and functions $f_1, \dots, f_K : S \times A \rightarrow \mathbb{R}$ that satisfy $\max_{k=1, \dots, K} \|k f_k - T^\pi f_k\|_{L_2(d_{p_0}^\pi)} \leq \eta$, and $f_0(s, a) = 0$. Then, for all $k = 1, \dots, K$, we have $\|k f_k - T^\pi f_k\|_{L_2(d_{p_0}^\pi)} \leq \frac{4}{1-\gamma} \eta + \gamma^{k/2}$.*

Proof. For any $k = 1, \dots, K$,

$$\begin{aligned} \|k f_k - T^\pi f_k\|_{L_2(d_{p_0}^\pi)} &\leq \|k f_k - T^\pi f_k\|_{L_2(d_{p_0}^\pi)} + \|k T^\pi f_k - T^\pi f_k\|_{L_2(d_{p_0}^\pi)} \\ &\leq \eta + \gamma (E_{s, a \sim d_{p_0}^\pi} [(E_{s^\theta, a^\theta \sim P(s, a) \sim \pi} [f_k(s^\theta, a^\theta) - f_k(s^\theta, a^\theta)])^2])^{1/2} \\ &\leq \eta + \gamma (E_{s, a \sim d_{p_0}^\pi, s^\theta, a^\theta \sim P(s, a) \sim \pi} [(f_k(s^\theta, a^\theta) - f_k(s^\theta, a^\theta))^2])^{1/2} \end{aligned}$$

Since $d_{p_0}^\pi(s, a) = \gamma \mathbb{E}_{s, a} d_{p_0}^\pi P(s|s, a)\pi(a|s) + (1 - \gamma)p_0(s)\pi(a|s)$, and the quantity inside the expectation is a square, and thus non-negative,

$$\begin{aligned} & \eta + \gamma(\gamma^{-1} \mathbb{E}_{s, a} d_{p_0}^\pi [(f_{k-1}(s, a) - f_k(s, a))^2])^{1/2} \\ = & \eta + \rho_{\bar{\gamma}} k_{f_{k-1} - f_k} k_{L_2(d_{p_0}^\pi)} \\ & \eta + \rho_{\bar{\gamma}} \left(\eta + k_{f_{k-1} - T^\pi f_{k-1}} k_{L_2(d_{p_0}^\pi)} \right). \end{aligned}$$

Unrolling the recursion and using $k_{f_0 - T^\pi f_0} k_{L_2(d_{p_0}^\pi)} = 1$, we have

$$\begin{aligned} k_{f_K - T^\pi f_K} k_{L_2(d_{p_0}^\pi)} &= \eta + \rho_{\bar{\gamma}} (\eta + \rho_{\bar{\gamma}} (\eta + \dots \rho_{\bar{\gamma}} (\eta + 1) \dots)) \\ &= \frac{1}{1 - \rho_{\bar{\gamma}}^K} \rho_{\bar{\gamma}}^K \eta + \gamma^{K/2}, \end{aligned}$$

which gives the claim since $\frac{1}{1 - \rho_{\bar{\gamma}}^K} = 2(1 - \gamma)^{-1}$ and $1 - \rho_{\bar{\gamma}}^K = 1$. □

The following lemma is a “fast rates”-like result for norms. It shows that the norm induced by the empirical covariance matrix $\hat{\Sigma}$ can be bounded by the norm induced by two times the population covariance matrix Σ , up to some $\tilde{O}(N^{-1/2})$ terms.

Lemma C.3 (Fast Rates for $\|\cdot\|_{\hat{\Sigma}}$ -norm). *For any $\delta \in (0, 1)$, we have with probability at least $1 - \delta$, for any $x \in B_W$, we have*

$$\|x\|_{\hat{\Sigma}} \leq 2\|x\|_{\Sigma} + 5W \sqrt{\frac{d \log(N/\delta)}{N}},$$

and

$$\|x\|_{\Sigma} \leq 2\|x\|_{\hat{\Sigma}} + 5W \sqrt{\frac{d \log(N/\delta)}{N}}.$$

Proof. First, fix any $x \in B_W$. Since $(x^\top \phi(s, a))^2 \leq W^2$ almost surely, we have $\mathbb{E} [(x^\top \phi(s, a))^4] \leq W^2 \mathbb{E} [(x^\top \phi(s, a))^2]$. So by Lemma B.1, w.p. at least $1 - \delta$,

$$\left| \frac{1}{N} \sum_{i=1}^N (x^\top \phi(s_i, a_i))^2 - \mathbb{E} [(x^\top \phi(s, a))^2] \right| \leq \frac{1}{2} \mathbb{E} [(x^\top \phi(s, a))^2] + \frac{2W^2 \log(2/\delta)}{N}.$$

In particular, this means $\|x\|_{\hat{\Sigma}} \leq \frac{3}{2}\|x\|_{\Sigma} + \frac{2W^2 \log(2/\delta)}{N}$ and $\|x\|_{\Sigma} \leq 2\|x\|_{\hat{\Sigma}} + \frac{4W^2 \log(2/\delta)}{N}$. Now union bound over a W/\sqrt{N} -net of B_W , which can be done with $(1 + 2/\sqrt{N})^d$ elements. The approximation error from this cover is

$$\begin{aligned} \|x\|_{\hat{\Sigma}} &\leq \|n(x)\|_{\hat{\Sigma}} + \|n(x) - x\|_{\hat{\Sigma}} \\ &\leq \sqrt{\frac{3}{2}} \|n(x)\|_{\Sigma} + \sqrt{\frac{2W^2 \log(2(1 + 2/\sqrt{N})^d/\delta)}{N}} + \frac{W}{\sqrt{N}} \\ &\leq \sqrt{\frac{3}{2}} \|x\|_{\Sigma} + \sqrt{\frac{3}{2}} \|n(x) - x\|_{\Sigma} + \sqrt{\frac{2W^2 \log(2(1 + 2/\sqrt{N})^d/\delta)}{N}} + \frac{W}{\sqrt{N}} \\ &\leq \sqrt{\frac{3}{2}} \|x\|_{\Sigma} + 4W \sqrt{\frac{d \log(N/\delta)}{N}}, \end{aligned}$$

where $n(x)$ is the closest element in the net to x . Similarly,

$$\begin{aligned} \|x - n(x)\| &\leq \|x - n(x)\| + \|n(x) - x\| \\ &\leq 2\|x - n(x)\| + \sqrt{\frac{4W^2 \log(2(1 + 2^{\frac{D}{N}})/\delta))}{N}} + \frac{W}{N} \\ &\leq 2\|x - n(x)\| + 2\|n(x) - x\| + \sqrt{\frac{4W^2 \log(2(1 + 2^{\frac{D}{N}})/\delta))}{N}} + \frac{W}{N} \\ &\leq 2\|x - n(x)\| + 5W\sqrt{\frac{d \log(N/\delta)}{N}}. \end{aligned}$$

□

Now define the following notation for analyzing Least Squares Policy Evaluation (LSPE). For every target vector $k\vartheta k_2$ W , define

$$\begin{aligned} y^\vartheta &= r + \gamma \vartheta^\top \phi(s^\vartheta, \pi), \\ \theta_\vartheta &= \arg \min_{k\theta k_2 \leq W} \mathbb{E}_{s,a} \mathbb{E}_{\nu, s^\vartheta \sim P(s,a)} [(y^\vartheta - \theta^\top \phi(s, a))^2] := \ell(\theta, \vartheta), \\ y_i^\vartheta &= r_i + \gamma \vartheta^\top \phi(s_i^\vartheta, \pi), \\ \widehat{\theta}_\vartheta &= \arg \min_{k\theta k_2 \leq W} \frac{1}{N} \sum_{i=1}^N (y_i^\vartheta - \theta^\top \phi(s_i, a_i))^2 := \widehat{\ell}(\theta, \vartheta). \end{aligned}$$

The following lemma are useful facts about the optimal θ_ϑ and $\widehat{\theta}_\vartheta$.

Lemma C.4. For any ϑ_1, ϑ_2 , we have

$$\begin{aligned} k\theta_{\vartheta_1} - \theta_{\vartheta_2} k_2 &\leq \sqrt{2\gamma W k\vartheta_1 - \vartheta_2 k_2}, \\ k\widehat{\theta}_{\vartheta_1} - \widehat{\theta}_{\vartheta_2} k_2 &\leq \sqrt{2\gamma W k\vartheta_1 - \vartheta_2 k_2}. \end{aligned}$$

For any θ, ϑ , we have

$$\ell(\theta, \vartheta) - \ell(\theta_\vartheta, \vartheta) \leq k\theta - \theta_\vartheta k_2^2.$$

Proof. First, recall that θ_ϑ minimizes $f(\theta) = \mathbb{E}[(y^\vartheta - \theta^\top \phi(s, a))^2]$, which has the Jacobian $\nabla f(\theta) = 2\mathbb{E}[(\theta^\top \phi(s, a) - y^\vartheta)\phi(s, a)]$. Since θ_ϑ is optimal over B_W , the necessary optimality condition is that $\nabla f(\theta_\vartheta) \succeq N_{B_W}(\theta_\vartheta)$, the normal cone of B_W at θ_ϑ , i.e. for any $\theta \succeq B_W$, we have $\langle \nabla f(\theta_\vartheta), \theta - \theta_\vartheta \rangle \geq 0$. In particular, we have

$$\begin{aligned} \langle \nabla f(\theta_\vartheta), \theta_{\vartheta_1} - \theta_\vartheta \rangle &\geq 0 \\ \langle \nabla f(\theta_\vartheta), \theta_{\vartheta_2} - \theta_\vartheta \rangle &\geq 0 \end{aligned}$$

adding the two we get

$$\begin{aligned} k\theta_{\vartheta_1} - \theta_{\vartheta_2} k_2 &= \langle \nabla f(\theta_\vartheta), \theta_{\vartheta_1} - \theta_{\vartheta_2} \rangle \leq \langle \nabla f(\theta_\vartheta), \theta_{\vartheta_1} - \theta_\vartheta \rangle + \langle \nabla f(\theta_\vartheta), \theta_\vartheta - \theta_{\vartheta_2} \rangle \\ &\leq \langle \nabla f(\theta_\vartheta), \theta_{\vartheta_1} - \theta_\vartheta \rangle + \langle \nabla f(\theta_\vartheta), \theta_\vartheta - \theta_{\vartheta_2} \rangle \\ &= \gamma \langle \nabla f(\theta_\vartheta), \phi(s, a)\phi(s^\vartheta, \pi)^\top \rangle (\vartheta_1 - \vartheta_2) \leq \gamma \langle \nabla f(\theta_\vartheta), \phi(s, a)\phi(s^\vartheta, \pi)^\top \rangle k_2 k\vartheta_1 - \vartheta_2 k_2 \\ &\leq \gamma \langle \nabla f(\theta_\vartheta), \phi(s, a)\phi(s^\vartheta, \pi)^\top \rangle k_2 k\vartheta_1 - \vartheta_2 k_2 \\ &\leq 2\gamma W k\vartheta_1 - \vartheta_2 k_2. \end{aligned}$$

The claim with $\widehat{\theta}_{\vartheta_1}, \widehat{\theta}_{\vartheta_2}$ follows by the same argument.

For the second claim, we first apply the Parallelogram law, followed by the first-order optimality of θ_ϑ ,

$$\begin{aligned} \ell(\theta, \vartheta) &= \ell(\theta_\vartheta, \vartheta) + \mathbb{E}_\nu((\theta_\vartheta - \theta)^\top \phi(s, a))^2 + 2\mathbb{E}_\nu[(y^\vartheta - \theta^\top \phi(s, a))\phi(s, a)^\top (\theta_\vartheta - \theta)] \\ &= \ell(\theta_\vartheta, \vartheta) + \mathbb{E}_\nu((\theta_\vartheta - \theta)^\top \phi(s, a))^2 + 0 \\ &= \ell(\theta_\vartheta, \vartheta) + k\theta - \theta_\vartheta k^2. \end{aligned}$$

□

Now we show our key lemma about the concentration of least squares, uniformly over all targets generated by $\vartheta \in B_W$.

Lemma C.5 (Concentration for Least Squares). *Let $F = \{f(s, a) - \theta^\top \phi(s, a) : k\theta k \leq W\}$, with $k\phi(s, a)k \leq 1$. Then, for any $\delta \in (0, 1)$ w.p. at least $1 - \delta$, we have*

$$\sup_{\vartheta \in B_W} k\widehat{\theta}_\vartheta - \theta_\vartheta k \leq 120d(1 + W) \frac{\log(N) \sqrt{\log(10/\delta)}}{N}.$$

Proof. By Lemma C.3, we have that w.p. at least $1 - \delta$, we can bound the random k -norm by k , and vice versa, simultaneously over all vectors in a ball,

$$E = \{f(s, a) - \theta^\top \phi(s, a) : k\theta k \leq t \text{ and } k\phi(s, a)k \leq t\},$$

provided $t \geq 5W \sqrt{\frac{d \log(N/\delta)}{N}}$. For the remainder of this proof, we *condition* on this high-probability event. That is, any probability and expectations will be *implicitly conditioned* on E .

First, we'll show that for an arbitrary and fixed $\vartheta \in B_W$, w.p. at least $1 - \delta$, we have

$$k\widehat{\theta}_\vartheta - \theta_\vartheta k \leq t.$$

To do so, we will bound the probability of the complement, which in turn can be simplified by the following chain of arguments,

$$k\widehat{\theta}_\vartheta - \theta_\vartheta k \leq t$$

By $\ell(\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + k\theta - \theta_\vartheta k^2$ from Lemma C.4,

$$\begin{aligned} &\Rightarrow \ell(\widehat{\theta}_\vartheta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + t^2 \\ &\Rightarrow \exists k\theta k \leq W : \ell(\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + t^2, \widehat{\ell}(\theta, \vartheta) = \widehat{\ell}(\theta_\vartheta, \vartheta) + 0 \end{aligned}$$

By convexity and continuity of $\ell(\cdot, \vartheta)$, we can make this strict equality. Indeed, given this, by Intermediate Value Theorem, there exists $\lambda \in [0, 1]$ such that $\theta^\theta = (1 - \lambda)\theta + \lambda\theta_\vartheta$ has $\ell(\theta^\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + \nu$. Then by convexity, $\widehat{\ell}(\theta^\theta, \vartheta) = \widehat{\ell}(\theta_\vartheta, \vartheta) + (1 - \lambda)\widehat{\ell}(\theta, \vartheta) + \lambda\widehat{\ell}(\theta_\vartheta, \vartheta) = 0$.

$$\begin{aligned} &\Rightarrow \exists k\theta k \leq W : \ell(\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + t^2, \widehat{\ell}(\theta, \vartheta) = \widehat{\ell}(\theta_\vartheta, \vartheta) + 0 \\ &\Rightarrow \exists k\theta k \leq W : k\theta - \theta_\vartheta k \leq t, (\ell(\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + t^2, \widehat{\ell}(\theta, \vartheta) = \widehat{\ell}(\theta_\vartheta, \vartheta) + 0) \end{aligned}$$

By conditioning on E ,

$$\Rightarrow \exists k\theta k \leq W : k\theta - \theta_\vartheta k \leq 3t, (\ell(\theta, \vartheta) = \ell(\theta_\vartheta, \vartheta) + t^2, \widehat{\ell}(\theta, \vartheta) = \widehat{\ell}(\theta_\vartheta, \vartheta) + 0)$$

Hence, we now focus on bounding

$$\mathbb{P}\left(\sup_{\theta \in E} \left| \sum_{i=1}^N \zeta_i(\theta, \vartheta) - \mathbb{E}_\nu[\zeta_i(\theta, \vartheta)] \right| \geq Nt^2\right) \leq \delta,$$

where we define

$$\begin{aligned}\zeta_i(\theta, \vartheta) &= (y_i^\vartheta - \theta^\top \phi(s_i, a_i))^2, \\ &= \|\theta - \theta_\vartheta\|_W^2, \quad \|\theta - \theta_\vartheta\|_W \leq 3tg.\end{aligned}\quad (8)$$

Since $\zeta_i(\theta, \vartheta) = \lambda_i(\|\theta - \theta_\vartheta\|_W)$ where $\lambda_i(x) = x(2y_i^\vartheta - (\theta + \theta_\vartheta)^\top \phi(s_i, a_i))$ is $2(1+W)$ -Lipschitz and $\lambda_i(0) = 0$, the contraction lemma (Lemma B.3) tells us that it suffices to consider a simpler class. Define

$$\omega_i(\theta, \vartheta) = (\theta - \theta_\vartheta)^\top \phi(s_i, a_i)$$

and $\omega(\theta, \vartheta) = (\omega_i(\theta, \vartheta))_{i=1}^N$. By Chaining (Lemma B.6), we have

$$\begin{aligned}\mathbb{E} \left[\left(\frac{1}{J} \sup_{\theta, \vartheta} \sum_{j=1}^J \omega(\theta, \vartheta) \right)^2 \right] &\leq 1 \\ \text{where } J &= \inf_{\alpha > 0} \left\{ 3\alpha \frac{\rho}{N} + 9 \int_{\alpha}^b \sqrt{\log(D(\delta/2, \omega(\cdot)))} d\delta \right\}, \\ \omega(\cdot) &= \mathcal{F}\omega(\theta, \vartheta) : \|\theta - \theta_\vartheta\|_W \leq g,\end{aligned}$$

$D(\delta, F)$ is the Euclidean packing number of $F \subset \mathbb{R}^N$, and $b = \sup_{\theta, \vartheta} \|\omega(\theta, \vartheta)\|_2$ is the envelope.

Now we'll bound the truncated entropy integral, J . First notice that $b \leq 3t \frac{\rho}{N}$, based on the definition of ζ_i (Equation (8)), which localizes in $\|\theta - \theta_\vartheta\|_W$,

$$\frac{b}{N} = \sup_{\theta, \vartheta} \sqrt{\frac{1}{N} \sum_{i=1}^N (\theta - \theta_\vartheta)^\top \phi(s_i, a_i) \phi(s_i, a_i) (\theta - \theta_\vartheta)} = \sup_{\theta, \vartheta} \|\theta - \theta_\vartheta\|_W \leq 3t.$$

Now, we bound the packing number $D(\cdot, \omega(\cdot))$. Let $\theta_1, \theta_2 \in B_W$ be arbitrary,

$$\frac{\|\omega(\theta_1, \vartheta) - \omega(\theta_2, \vartheta)\|_2}{N} = \|\theta_1 - \theta_2\|_W \sqrt{\sigma_{\max}(\widehat{\Sigma})} \leq \|\theta_1 - \theta_2\|_W \sqrt{\sigma_{\max}(\widehat{\Sigma})}.$$

So for any ε , we can construct an ε -cover by setting $\|\theta_1 - \theta_2\|_W \leq \varepsilon / \sqrt{\sigma_{\max}(\widehat{\Sigma})}$, which requires $(W(1 + 2 \frac{\rho}{N}/\varepsilon))^d$ points. Let $N(\varepsilon, F)$ denote the Euclidean covering number of $F \subset \mathbb{R}^N$. Then,

$$\log D(\varepsilon, \omega(\cdot)) \leq \log N(\varepsilon/2, \omega(\cdot)) \leq d \log W(1 + 4 \frac{\rho}{N}/\varepsilon)$$

So,

$$\begin{aligned}J &\leq \int_0^{3t \frac{\rho}{N}} \sqrt{\log(D(\varepsilon/2, \omega(\cdot)))} d\varepsilon \\ &\leq \int_0^{3t \frac{\rho}{N}} \sqrt{d \log(W(1 + 8 \frac{\rho}{N}/\varepsilon))} d\varepsilon \\ &\leq 3t \frac{\rho}{N} (\sqrt{\log W} + \int_0^1 \sqrt{\log(1 + 3/(\varepsilon t))} d\varepsilon) \\ &\leq 3t \frac{\rho}{N} (\sqrt{\log W} + \sqrt{\log(4/t)}),\end{aligned}$$

since $\int_0^1 \sqrt{\log(1 + c/\varepsilon)} d\varepsilon \leq \sqrt{\log(1 + c)}$ for any $c > 0$, and assuming $t \geq 1$.

Now we put everything together. Let c denote a positive constant,

$$\mathbb{P} \left(\sup_{\theta, \vartheta} \left| \sum_{i=1}^N \zeta_i(\theta, \vartheta) - \mathbb{E} \zeta_i(\theta, \vartheta) \right| \geq Nt^2 \right)$$

Since ϵ is increasing,

$$= \mathbb{P} \left(\left(c \sup \left| \sum_{i=1}^N \zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta) - \mathbb{E}_\nu [\zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta)] \right| \right) > (cNt^2) \right)$$

By Markov's inequality,

$$\mathbb{E} \left[\frac{\left(c \sup \left| \sum_{i=1}^N \zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta) - \mathbb{E}_\nu [\zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta)] \right| \right)}{(cNt^2)} \right]$$

By Symmetrization (Lemma B.2),

$$\mathbb{E} \left[\mathbb{E} \left[\frac{(2c \sup_{\theta, \vartheta} \left| \sum_{i=1}^N \zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta) \right|)}{(cNt^2)} \right) \right]$$

By Contraction (Lemma B.3) and that $\zeta_i(\theta, \vartheta) - \zeta_i(\theta_\vartheta, \vartheta) = \lambda_i(\omega_i(\theta, \vartheta))$ where λ_i is $L = 2(1 + W)$ Lipschitz,

$$\frac{3}{2} \mathbb{E} \left[\mathbb{E} \left[\frac{(4cL \sup_{\theta, \vartheta} \left| \sum_{i=1}^N \omega_i(\theta, \vartheta) \right|)}{(cNt^2)} \right) \right]$$

Setting $c = \frac{1}{4LJ}$ and applying Chaining (Lemma B.6), the numerator is bounded by 1,

$$8 \exp \left(- \left(\frac{Nt^2}{4LJ} \right)^2 \right)$$

Applying upper bound on J ,

$$8 \exp \left(- \left(\frac{Nt^2}{8(1+W) \cdot 3t \cdot \frac{1}{dN} (\frac{1}{\log W} + \sqrt{\log(4/t)})} \right)^2 \right)$$

$$8 \exp \left(- \frac{Nt^2}{24^2(1+W)^2 d(\log W + \log(4/t))} \right)$$

Now, we set $t = 24(1+W) \sqrt{\frac{d \log(N) \log(1/\delta)}{N}}$,

$$8 \exp \left(- \frac{\log(N) \log(1/\delta)}{\log(N/d \log(N))} \right)$$

Since $\log(N) \geq \log(N/d \log(N))$,

$$8\delta.$$

Hence, we have shown that for an arbitrary and fixed $\vartheta \in B_W$, w.p. $1 - 9\delta$, we have

$$\left| \widehat{\theta}_\vartheta - \theta_\vartheta \right| < t = 24(1+W) \sqrt{\frac{d \log(N) \log(1/\delta)}{N}}.$$

We finally apply a union bound over $\vartheta \in B_W$. Consider a W/N -cover of $\vartheta \in B_W$, which requires $(1 + 2N)^d$ points. Then, for any $\vartheta \in B_W$, we have

$$\left| \widehat{\theta}_\vartheta - \theta_\vartheta \right| \leq \left| \widehat{\theta}_\vartheta - \widehat{\theta}_{n(\vartheta)} \right| + \left| \widehat{\theta}_{n(\vartheta)} - \theta_{n(\vartheta)} \right| + \left| \theta_{n(\vartheta)} - \theta_\vartheta \right|$$

$$\leq (2\widehat{\theta}_\vartheta - \widehat{\theta}_{n(\vartheta)} + t) + t + \sqrt{2\gamma W k n(\vartheta)} \cdot \vartheta$$

$$\leq 2t + 3\sqrt{2\gamma W k n(\vartheta)} \cdot \vartheta$$

$$\leq 2t + 3\sqrt{2\gamma W^2/N}$$

$$\leq 5t.$$

Thus, we have shown that w.p. $1 - \delta$,

$$\sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k < 120d(1+W) \frac{\log(N)}{\beta} \sqrt{\frac{\log(10/\delta)}{N}}.$$

□

A nice corollary is that when β provides full coverage, i.e. $\beta \geq \beta_I$ for some positive β , then we can bound

$$\sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k \leq \sigma_{\min}(\cdot)^{-1/2} \sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k \leq \beta^{-1/2} \sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k.$$

C.1. Main LSPE theorem

We now prove our LSPE sample complexity guarantee Theorem 3.3.

Theorem 3.3 (Sample Complexity of LSPE). *Assume feature ϕ satisfies approximate Linear BC with parameter ε_ν . For any $\delta \in (0, 0.1)$, w.p. at least $1 - \delta$, we have for any initial state distribution p_0*

$$\begin{aligned} & \left| V_{p_0}^{\pi_e} - \mathbb{E}_{s \sim p_0} [\hat{f}_K(s, \pi_e)] \right| \leq \frac{\gamma^{K/2}}{1 - \gamma} + \frac{4 \sqrt{\left\| \frac{dd_{p_0}^{\pi_e}}{d\nu} \right\|_1}}{(1 - \gamma)^2} \varepsilon_\nu \\ & + \frac{480 \sqrt{\kappa(p_0)} (1+W) d \log(N) \sqrt{\log(10/\delta)}}{(1 - \gamma)^2 \beta \bar{N}}, \end{aligned}$$

where \hat{f}_K is the output of Algorithm 1.

Proof of Theorem 3.3. By Lemmas C.1 and C.2,

$$\left| V_{p_0}^{\pi_e} - \mathbb{E}_{s \sim p_0} [\hat{f}_k(s, \pi)] \right| \leq \frac{4}{(1 - \gamma)^2} \max_{k=1,2,\dots} k \hat{f}_k - T^\pi \hat{f}_k - 1 k_{L_2(d_{p_0}^\pi)} + \frac{\gamma^{k/2}}{1 - \gamma}.$$

Next, we bound the maximum regression error. Consider any initial state distribution p_0 , then we have

$$\begin{aligned} & \max_{k=1,2,\dots} k \hat{f}_k - T^\pi \hat{f}_k - 1 k_{L_2(d_{p_0}^\pi)} \leq \sup_{k \in B_W} k_{\hat{\theta}_\vartheta} \phi - T^\pi(\vartheta^\dagger \phi) k_{L_2(d_{p_0}^\pi)} \\ & \leq \sup_{k \in B_W} k_{\hat{\theta}_\vartheta} \phi - \theta_\vartheta^\dagger \phi k_{L_2(d_{p_0}^\pi)} + \sup_{k \in B_W} k_{\theta_\vartheta^\dagger} \phi - T^\pi(\vartheta^\dagger \phi) k_{L_2(d_{p_0}^\pi)} \\ & \leq \sqrt{\kappa(p_0)} \sup_{k \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k + \sqrt{\left\| \frac{dd_{p_0}^\pi}{d\nu} \right\|_1} \sup_{k \in B_W} k_{\theta_\vartheta^\dagger} \phi - T^\pi(\vartheta^\dagger \phi) k_{L_2(\nu)} \\ & \leq \sqrt{\kappa(p_0)} \sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k + \sqrt{\left\| \frac{dd_{p_0}^\pi}{d\nu} \right\|_1} \varepsilon_\nu. \end{aligned}$$

The quantity $\sup_{\vartheta \in B_W} k_{\hat{\theta}_\vartheta} - \theta_\vartheta k$ can be directly bounded by Lemma C.5 w.p. at least $1 - \delta$. Thus, we have shown the desired result: for any initial state distribution p_0 ,

$$\left| V_{p_0}^{\pi_e} - \mathbb{E}_{s, a \sim p_0, \pi} [\hat{f}_k(s, a)] \right| \leq \frac{4}{(1 - \gamma)^2} \left(\sqrt{\kappa(p_0)} 120d(1+W) \frac{\log(N)}{\beta} \sqrt{\frac{\log(10/\delta)}{N}} + \sqrt{\left\| \frac{dd_{p_0}^\pi}{d\nu} \right\|_1} \right) \varepsilon_\nu + \frac{\gamma^{k/2}}{1 - \gamma}.$$

□

D. Proofs for Linear BC Equivalence

Proposition 4.2. Consider a feature ϕ with full rank covariance $\Sigma(\phi)$. Given any $W > 0$, the feature ϕ being linear BC (under B_W) implies that there exist $(\rho, M) \succeq B_W \in \mathbb{R}^d \times \mathbb{R}^{d \times d}$ with $kMk_2 \leq \sqrt{1 - \frac{k\rho k_2^2}{W^2}}$, and

$$\mathbb{E}_\nu \left\| \begin{bmatrix} M \\ \rho^\top \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \mathbb{E}_{s^0 \sim P(s,a)} \phi(s^0, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 = 0. \quad (2)$$

On the other hand, if there exists $(\rho, M) \succeq B_W \in \mathbb{R}^d \times \mathbb{R}^{d \times d}$ with $kMk_2 < 1$ such that the above equality holds, then ϕ must satisfy exact linear BC with $W = \frac{k\rho k_2^2}{1 - kMk_2}$.

Proof. (\Leftarrow) Suppose that $kMk_2 < 1$ and $\rho \succeq B_W$ satisfy,

$$\mathbb{E}_\nu \left\| \begin{bmatrix} M \\ \rho^\top \end{bmatrix} \phi(s, a) - \begin{bmatrix} \gamma \mathbb{E}_{s^0 \sim P(s,a)} \phi(s^0, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 = 0.$$

Then, for any $w_1 \succeq B_W$, setting $w_2 = \rho + M^\top w_1$ satisfies,

$$\begin{aligned} kw_2^\top \phi - T^{\pi_e}(w_1^\top \phi) k_\nu^2 &= \mathbb{E}_\nu \left\| w_2^\top \phi(s, a) - r(s, a) - \gamma \mathbb{E}_{s^0 \sim P(s,a)} [w_1^\top \phi(s^0, \pi_e)] \right\|_2^2 \\ &= \mathbb{E}_\nu \left\| \rho^\top \phi(s, a) - r(s, a) + w_1^\top M \phi(s, a) - \gamma \mathbb{E}_{s^0 \sim P(s,a)} [w_1^\top \phi(s^0, \pi_e)] \right\|_2^2 \\ &= \mathbb{E}_\nu \left\| w_1^\top (M \phi(s, a) - \gamma \mathbb{E}_{s^0 \sim P(s,a)} [\phi(s^0, \pi_e)]) \right\|_2^2 \\ &= 0. \end{aligned}$$

Also, we have $k\rho k_2 \leq k\rho k_2 + kMk_2 k w_1 k_2 \leq W$ since $W = \frac{k\rho k_2^2}{1 - kMk_2}$. Thus, ϕ satisfies exact Linear BC.

(\Rightarrow) Suppose ϕ satisfies exact Linear BC, that is

$$\max_{w_1 \succeq B_W} \min_{w_2 \succeq B_W} kw_2^\top \phi - T^\pi(w_1^\top \phi) k_\nu^2 = 0.$$

To see that there exists $\rho \succeq B_W$ that linearizes the reward w.r.t ϕ under ν , set $w_1 = 0$, and we have:

$$\min_{w_2 \succeq B_W} \mathbb{E}_{s,a \sim \nu} \|w_2^\top \phi(s, a) - r(s, a)\|_2^2 = 0.$$

Let ρ to be the minimizer of the above objective.

Now we need to show that there exists a $M \succeq \mathbb{R}^{d \times d}$ with $kMk_2 < \sqrt{1 - \frac{k\rho k_2^2}{W^2}}$ that satisfies

$$\mathbb{E}_{s,a \sim \nu} \left\| M \phi(s, a) - \gamma \mathbb{E}_{s^0 \sim P(s,a)} \phi(s^0, \pi_e) \right\|_2^2 = 0.$$

To extract the i -th row of M , plug in $w_i = W e_i$ (note that $w_i \succeq B_W$). By exact Linear BC, we know that there exists a vector $v_i \succeq B_W$, such that:

$$\|v_i^\top \phi(s, a) - \rho^\top \phi(s, a) - \gamma W \mathbb{E}_{s^0 \sim P(s,a)} e_i^\top \phi(s^0, \pi_e)\|_\nu = 0.$$

Repeating this for every $i \in [d]$, we can construct M as follows,

$$M = \frac{1}{W} \begin{bmatrix} (v_1 \ \rho^\top)^\top \\ \vdots \\ (v_d \ \rho^\top)^\top \end{bmatrix},$$

which satisfies,

$$\begin{aligned}
 & \mathbb{E}_{s,a} \nu \left\| M\phi(s,a) - \gamma \mathbb{E}_{s^0} P(s,a)\phi(s^0, \pi_e) \right\|_2^2 \\
 &= \sum_{i=1}^d \mathbb{E}_{s,a} \nu \left\| e_i^\top (M\phi(s,a) - \gamma \mathbb{E}_{s^0} P(s,a)\phi(s^0, \pi_e)) \right\|_2^2 \\
 &= \sum_{i=1}^d \mathbb{E}_{s,a} \nu \left\| \frac{1}{W} (v_i - \rho)^\top \phi(s,a) - \gamma \mathbb{E}_{s^0} P(s,a) e_i^\top \phi(s^0, \pi_e) \right\|_2^2 = 0.
 \end{aligned}$$

Hence, we have $\left\| M\phi(s,a) - \gamma \mathbb{E}_{s^0} P(s,a)\phi(s^0, \pi_e) \right\|_2^2 = 0$.

Finally, we must show that $kMk_2 < \sqrt{1 - \frac{k\rho k_2^2}{W^2}}$. First we show that $kMk_2 \leq 1$. For any $w_1 \in B_W$, by exact linear BC, there exists $w_2 \in B_W$ s.t. $\left\| w_2^\top \phi(s,a) - r(s,a) - \gamma \mathbb{E}_{s^0} P(s,a) [w_1^\top \phi(s^0, \pi_e)] \right\|_2^2 = 0$, and by the construction of M , satisfies $k(w_2 - \rho - M^\top w_1)^\top \phi(s,a) k_2^2 = k w_2 - \rho - M^\top w_1 k_2^2 = 0$. Since ϕ is positive definite, we have that $w_2 = \rho + M^\top w_1$ is the unique choice of w_2 , which by exact linear BC is in B_W . Hence, we've shown that for any $w_1 \in B_W$, we also have that $\rho + M^\top w_1 \in B_W$. Now take w_1 and $-w_1$, subtracting the two expressions yields that $2M^\top w_1 \in B_W$. Since this is true for arbitrary w_1 , taking supremum over $w_1 \in B_W$ shows that $kMk_2 \leq 1$.

Now we show that the inequality must be strict. Consider the singular value decomposition: $M = \sum_{i=1}^d \sigma_i u_i v_i^\top$ where $\{u_i\}_{i \in [d]}$ and $\{v_i\}_{i \in [d]}$ are each an orthonormal basis of \mathbb{R}^d , and σ_i is the i -th largest singular value. Without loss of generality, suppose $\rho^\top u_1 > 0$, since we can always flip the sign of u_1 . If we pick $x = Wv_1 \in B_W$, we have $Mx = W\sigma_1 u_1$. By the argument in the previous paragraph, since $x \in B_W$, we have $\rho + Mx \in B_W$, implying that

$$W^2 - kMx + \rho k_2^2 = kW\sigma_1 u_1 + (\rho^\top u_1)u_1 + (\rho - (\rho^\top u_1)u_1) k_2^2$$

Since u_1 and $(\rho - (\rho^\top u_1)u_1)$ are orthogonal, by Pythagoras, we have

$$\begin{aligned}
 &= \int W\sigma_1 + \rho^\top u_1 \int^2 + k(\rho - (\rho^\top u_1)u_1) k_2^2 \\
 &= (W\sigma_1)^2 + 2W\sigma_1 \rho^\top u_1 + (\rho^\top u_1)^2 + k(\rho - (\rho^\top u_1)u_1) k_2^2 \\
 &= (W\sigma_1)^2 + 2W\sigma_1 \rho^\top u_1 + k(\rho^\top u_1)u_1 k_2^2 + k(\rho - (\rho^\top u_1)u_1) k_2^2
 \end{aligned}$$

By Pythagoras,

$$= (W\sigma_1)^2 + 2W\sigma_1 \rho^\top u_1 + k\rho k_2^2.$$

Hence, we get the following inequality:

$$W^2 \sigma_1^2 + 2W(\rho^\top u_1)\sigma_1 + k\rho^2 k_2^2 - W^2 \leq 0.$$

Solve for σ_1 and using the fact that $\rho^\top u_1 > 0$, we have that,

$$\sigma_1 = \frac{\rho^\top u_1 + \sqrt{(\rho^\top u_1)^2 + (W^2 - k\rho k_2^2)}}{W} = \frac{\sqrt{W^2 - k\rho k_2^2}}{W} + \sqrt{1 - \frac{k\rho k_2^2}{W^2}}.$$

We finally show that $k\rho k_2 < W$ unless $M = 0$. We prove this by contradiction. Assume $k\rho k_2 = W$. Following the above argument, take any $w_1 \in B_W$, we must have $w_2 := \rho + M^\top w_1 \in B_W$. We discuss two cases.

First if $\rho \notin \text{range}(M^\top)$. In this case, we must have $k w_2 k_2^2 = k\rho k_2^2 + kM^\top w_1 k_2^2 = W^2 + kM^\top w_1 k_2^2 > W^2$, as long as M has non-zero entries. Thus, this case leads to contradiction.

Second, if $\rho \in \text{range}(M^\top)$. In this case, there must exist a vector $x \neq 0$ such that $M^\top x = \rho$. Consider $x := W \frac{x}{kx k_2} \in B_W$. We have $w_2 := \rho + M^\top x = \rho \left(1 + \frac{W}{kx k_2}\right)$, which means that $k w_2 k_2 > W$, which causes contradiction again.

So unless $M = 0$, which only happens when $\gamma = 0$ (i.e. horizon is 1), we have $k\rho k_2 < W$. \square

We now show an approximate version to the equivalence of Proposition 4.2. First, recall the bilevel loss from Equation (3). We use L_{lbc} and \widehat{L}_{lbc} to denote the population and empirical versions as follows,

$$L_{lbc}(\phi) = \min_{(\rho, M) \geq 2} E_{\nu} \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma E_{s^0} P(s, a) [\phi(s^0, \pi_e)] \\ r(s, a) \end{bmatrix} \right\|_2^2 \quad (9)$$

$$\widehat{L}_{lbc}(\phi) = \min_{(\rho, M) \geq 2} E_D \left\| \begin{bmatrix} M \\ \rho \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma \phi(s^0, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \min_{g \geq 2G} E_D kg(s, a) \quad \gamma \phi(s^0, \pi_e) k_2^2, \quad (10)$$

where $\geq = \{(\rho, M) \geq 2 B_W \quad \mathbb{R}^d \quad d : k\rho k \quad k\rho^* k, kM k_2 \quad kM^* k_2\}$ and (ρ^*, M^*) are corresponding to the linear BC ϕ^* .

Lemma D.1. *Suppose a feature $\widehat{\phi}$ satisfies $L_{lbc}(\widehat{\phi}) \leq \varepsilon^2$. Then, $\widehat{\phi}$ is $\varepsilon(1 + W)$ -approximately Linear BC, provided $W \leq \frac{k\rho^* k_2}{1 - kM^* k_2}$.*

Proof. Suppose $L_{lbc}(\widehat{\phi}) \leq \varepsilon^2$, so there exists \widehat{M} (s.t. $k\widehat{M} k_2 \quad kM^* k_2 < 1$) and $\widehat{\rho} \geq 2 B_W$ s.t.

$$E_{s, a} \nu \left\| \begin{bmatrix} \widehat{M} \\ \widehat{\rho} \end{bmatrix} \phi(s, a) \begin{bmatrix} \gamma E_{s^0} P(s, a) \phi(s^0, \pi_e) \\ r(s, a) \end{bmatrix} \right\|_2^2 \leq \varepsilon^2$$

For any $w_1 \geq 2 B_W$, we can take $w_2 = \widehat{\rho} + \widehat{M}^\top w_1$. Then, $k w_2 k_2 \leq k \widehat{\rho} k_2 + k \widehat{M} k_2 W \leq k \rho^* k_2 + k M^* k_2 W \leq W$ by our assumption on W . Hence,

$$\begin{aligned} & \max_{w_1 \geq 2 B_W} \min_{w_2 \geq 2 B_W} \left\| w_2^\top \phi(s, a) \quad r(s, a) \quad \gamma E_{s^0} P(s, a) [w_1^\top \phi(s^0, \pi_e)] \right\|_{\nu} \\ & \max_{w_1 \geq 2 B_W} \left\| (\widehat{\rho} + \widehat{M}^\top w_1)^\top \phi(s, a) \quad r(s, a) \quad \gamma E_{s^0} P(s, a) [w_1^\top \phi(s^0, \pi_e)] \right\|_{\nu} \\ & = \max_{w_1 \geq 2 B_W} \sqrt{E_{s, a} \nu \left[\left((\widehat{\rho} + \widehat{M}^\top w_1)^\top \phi(s, a) \quad r(s, a) \quad \gamma E_{s^0} P(s, a) [w_1^\top \phi(s^0, \pi_e)] \right)^2 \right]} \\ & \max_{w_1 \geq 2 B_W} \sqrt{E_{s, a} \nu \left[(\widehat{\rho}^\top \phi(s, a) \quad r(s, a))^2 \right]} + \sqrt{E_{s, a} \nu \left[\left(w_1^\top (\widehat{M} \phi(s, a) \quad \gamma E_{s^0} P(s, a) [\phi(s^0, \pi_e)]) \right)^2 \right]} \\ & \leq \varepsilon(1 + W), \end{aligned}$$

as desired. □

E. Proofs for Representation Learning

To simplify analysis, assume that the functions in G have bounded ℓ_2 norm, i.e. $\forall g \geq 2 G, s, a \geq 2 S \quad A, kg(s, a) k_2 \leq \gamma$. This is reasonable, and can always be achieved by clipping without loss of accuracy, since the target for $g(s, a)$ is $\gamma E_{s^0} P(s, a) [\phi(s^0, \pi_e)]$ and $k\phi(s, a) k_2 \leq 1$ for any $s, a \geq 2 S \quad A$.

E.1. Lemmas

Recall that $\lambda_k(A)$ denotes the k -th largest eigenvalue of a matrix A , i.e. $\lambda_1(A), \lambda_n(A)$ give the largest and smallest eigenvalues respectively.

Lemma E.1 (Weyl's Perturbation Theorem). *Let $A, B \geq 2 C^n \quad n$ be Hermitian matrices. Then*

$$\max_k |\lambda_k(A) - \lambda_k(B)| \leq kA - B k_2.$$

Proof. Please see (Bhatia, 2013, Corollary III.2.6). □

We extend this to be uniform over all $\phi \in \mathcal{F}$.

Lemma E.2 (Uniform spectrum concentration). *For any $\delta \in (0, 1)$, w.p. $1 - \delta$,*

$$\sup_{\phi \in \mathcal{F}, k \geq [d]} \left| \lambda_k(\phi) - \lambda_k(\widehat{\phi}) \right| \leq \frac{1}{N} \left(96\kappa(\bar{d}) + 4d + 4 \log^{1/2}(1/\delta) \right)$$

Proof. First observe that,

$$\begin{aligned} \sup_{\phi \in \mathcal{F}} \sup_k \left| \lambda_k(\phi) - \lambda_k(\widehat{\phi}) \right| &= \sup_{\phi \in \mathcal{F}} \left\| \phi - \widehat{\phi} \right\|_2 \\ &= \sup_{\phi \in \mathcal{F}, k \geq k_2 - 1} \left| x^\top (\phi - \widehat{\phi}) x \right| \\ &= \sup_{\phi \in \mathcal{F}, k \geq k_2 - 1} (E_\nu - E_D)(x^\top \phi(s, a))^2 \end{aligned}$$

Now we want to bound the Rademacher complexity of the class

$$F = \left\{ (s, a) \mapsto (x^\top \phi(s, a))^2, \phi \in \mathcal{F}, k \geq k_2 - 1 \right\}$$

First, to bound the envelope, we have $(x^\top \phi(s, a))^2 \leq 1$. To cover, consider any $\phi \in \mathcal{F}$ and $x \in \mathbb{R}^d$ s.t. $k \geq k_2 - 1$. Pick $\tilde{\phi}, \tilde{x}$ close to ϕ, x , so that

$$\begin{aligned} &\sqrt{\frac{1}{N} \sum_{i=1}^N \left((x^\top \phi(s_i, a_i))^2 - (\tilde{x}^\top \tilde{\phi}(s_i, a_i))^2 \right)^2} \\ &\leq 2 \sqrt{\frac{1}{N} \sum_{i=1}^N \left((x^\top \phi(s_i, a_i) - \tilde{x}^\top \tilde{\phi}(s_i, a_i)) \right)^2} \\ &\leq 2 \left(\sqrt{\frac{1}{N} \sum_{i=1}^N \left(x^\top (\phi(s_i, a_i) - \tilde{\phi}(s_i, a_i)) \right)^2} + \sqrt{\frac{1}{N} \sum_{i=1}^N \left((x - \tilde{x})^\top \tilde{\phi}(s_i, a_i) \right)^2} \right) \\ &\leq 2 \left(d(\phi, \tilde{\phi}) + k(x - \tilde{x})_{k_2} \right) \end{aligned}$$

So it suffices to take $d(\phi, \tilde{\phi}), k(x - \tilde{x})_{k_2} \leq t/4$ to t -cover F in $L_2(D)$. Note the $t/4$ -covering number for x in the unit ball is $(1 + 8/t)^d$. Thus, by Dudley's entropy bound ((5.48) of (Wainwright, 2019)),

$$\begin{aligned} R_N(F) &\leq \frac{24}{N} \int_0^1 \log^{1/2}(N(t/4, \cdot)) (1 + 8/t)^d dt \\ &\leq \frac{96}{N} \left(\kappa(\bar{d}) + 4 \sqrt{\bar{d}} \right) \end{aligned}$$

Thus, by Theorem 4.10 of (Wainwright, 2019), w.p. $1 - \delta$,

$$\begin{aligned} \sup_{\phi \in \mathcal{F}, k \geq k_2 - 1} (E_\nu - E_D)(x^\top \phi(s, a))^2 &\leq 2R_N(F) + \frac{4 \log^{1/2}(1/\delta)}{N} \\ &\leq \frac{1}{N} \left(96\kappa(\bar{d}) + 4 \sqrt{\bar{d}} + 4 \log^{1/2}(1/\delta) \right) \end{aligned}$$

□

We now prove the double sampling lemma, i.e. modified Bellman Residual Minimization (Chen & Jiang, 2019). This will help deal with the double sampling issue when transitions are stochastic. Recall that G is a function class of functions $g : X \mapsto \mathbb{R}^d$. Let ν be a distribution over $X \in \mathbb{R}^d$ and, for any $x \in X$, let $P(x)$ be a distribution over $Y \in \mathbb{R}^d$.

Lemma E.3 (Double Sampling). *Suppose $x \not\sim E_{y \sim P(x)} [y] \in G$. Then,*

$$E_{x \sim \nu} \left[\left\| x - E_{y \sim P(x)} [y] \right\|_2^2 \right] = E_{x \sim \nu, y \sim P(x)} \left[kx - yk_2^2 \right] \inf_{g \in G} E_{x \sim \nu, y \sim P(x)} \left[kg(x) - yk_2^2 \right]$$

Proof.

$$\begin{aligned} & E_{x \sim \nu, y \sim P(x)} \left[kx - yk_2^2 \right] - E_{x \sim \nu} \left[\left\| x - E_{y \sim P(x)} [y] \right\|_2^2 \right] \\ &= E_{x \sim \nu} \left[E_{y \sim P(x)} \left[kxk_2^2 - 2\langle x, y \rangle + kyk_2^2 \right] - \left(kxk_2^2 - 2\langle x, E_{y \sim P(x)} [y] \rangle + \left\| E_{y \sim P(x)} [y] \right\|_2^2 \right) \right] \\ &= E_{x \sim \nu} \left[E_{y \sim P(x)} \left[kyk_2^2 \right] - \left\| E_{y \sim P(x)} [y] \right\|_2^2 \right] \\ &= E_{x \sim \nu} \left[E_{y \sim P(x)} \left[\left\| y - E_{y \sim P(x)} [y] \right\|_2^2 \right] \right] \end{aligned}$$

where the last step uses the fact that $\left\| E_{y \sim P(x)} [y] \right\|_2^2 = E_{y \sim P(x)} [\langle y, E_{y \sim P(x)} [y] \rangle]$, and completing the square. Now, observe that, assuming $g^*(x) \in G$, we have that it is the minimizer of,

$$g^* \in \arg \min_{g \in G} E_{x \sim \nu, y \sim P(x)} \left[kg(x) - yk_2^2 \right] \quad (11)$$

which completes the proof. \square

E.2. Concentration lemmas

For any ϕ , define the optimal ρ, M, g for our losses as follows:

$$\begin{aligned} \rho_\phi &\in \arg \min_{(\rho, \cdot)} E_\nu \left[(\rho^\top \phi(s, a) - r(s, a))^2 \right] \\ M_\phi &\in \arg \min_{(\cdot, M)} E_\nu \left[kM\phi(s, a) - \gamma\phi(s^\theta, \pi_e)k_2^2 \right] \\ g_\phi &\in \arg \min_{g \in G} E_\nu \left[kg(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 \right] \end{aligned}$$

Similarly, define $\widehat{\rho}_\phi, \widehat{M}_\phi, \widehat{g}_\phi$ to be minimizers of the above losses when expectation is taken over the empirical distribution D , instead of the population distribution ν . Observe that the unconstrained minimization yields a closed form solution for g_ϕ as $g_\phi(s, a) = \gamma E_{s^\theta \sim P(s, a)} [\phi(s^\theta, \pi)]$ – Assumption 4.3 posits that G is rich enough to capture this.

The key property of our squared losses is that the second moment can be upper bounded by the expectation, which allows us to invoke the second part of the above Lemma B.1. We now combine this with covering to get uniform convergence results.

Lemma E.4. *For any $\delta \in (0, 1)$, w.p. at least $1 - \delta$, for any $\phi \in \mathcal{F}$, $\rho \in B_W$ and $M \in \mathbb{R}^{d \times d}$ with $kMk_2 \leq 1$, we have*

$$\begin{aligned} & \left| E_D \left[(\rho^\top \phi(s, a) - r(s, a))^2 \right] - E_\nu \left[(\rho^\top \phi(s, a) - r(s, a))^2 \right] \right| \leq \frac{1}{2} E_\nu \left[(\rho^\top \phi(s, a) - r(s, a))^2 \right] + \varepsilon_\rho, \\ & \left| E_D \left[kM\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 - kg_\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 \right] - E_\nu \left[kM\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 - kg_\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 \right] \right| \\ & \leq \frac{1}{2} E_\nu \left[kM\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 - kg_\phi(s, a) - \gamma\phi(s^\theta, \pi)k_2^2 \right] + \varepsilon_M, \end{aligned}$$

and, assuming realizability (Assumption 4.3), for every $g \in G$, we have

$$\begin{aligned} & \left| E_D \left[kg(s, a) - \gamma\phi^*(s^\theta, \pi)k_2^2 - kg_{\phi^*}(s, a) - \gamma\phi^*(s^\theta, \pi)k_2^2 \right] - E_\nu \left[kg(s, a) - g_{\phi^*}(s, a)k_2^2 \right] \right| \\ & \leq \frac{1}{2} E_\nu \left[kg(s, a) - g_{\phi^*}(s, a)k_2^2 \right] + \varepsilon_g, \end{aligned}$$

where $\varepsilon_\rho, \varepsilon_M, \varepsilon_g$ are defined below.

Finite function classes Assuming \mathcal{F} and G are finite, we have

$$\begin{aligned}\varepsilon_\rho & \frac{6d(1+W)^2 \log(4Wj/N/\delta)}{N} \\ \varepsilon_M & \frac{32d^2 \log(4j/N/\delta)}{N} \\ \varepsilon_g & \frac{20\gamma^2 \log(2jGj/\delta)}{N}.\end{aligned}$$

Proof for ε_ρ . For a fixed $\phi \in \mathcal{F}$, $\rho \in B_W$, apply Lemma B.1 to $X_i = (\rho^\top \phi(s_i, a_i) - r(s_i, a_i))^2$. The envelope is $jX_ij = (1+W)^2$ and the second moment is bounded $\mathbb{E}[X_i^2] \leq (1+W)^2 \mathbb{E}[X_i]$. So, the error from the lemma is $\frac{2(1+W)^2 \log(2/\delta)}{N}$. Now union bound over an ε -net of B_W . Since $|(\rho^\top \phi(s, a) - r(s, a))^2 - (\tilde{\rho}^\top \phi(s, a) - r(s, a))^2| = j(\rho - \tilde{\rho})^\top \phi(s, a)((\rho + \tilde{\rho})^\top \phi(s, a) + 2r(s, a))j \leq 2(1+W)k\tilde{\rho} - \rho k_2$, we consider a $\frac{1}{N}$ -net of B_W which requires $(1+2WN)^d$ points. The error from this ε -net approximation is at most $\frac{4(1+W)}{N}$. Finally, union bound over \mathcal{F} . \square

Proof for ε_M . For a fixed $\phi \in \mathcal{F}$ and $M \in \mathbb{R}^{d \times d}$ s.t. $kMk_2 < 1$, apply Lemma B.1 to $X_i = ka k_2^2 - kb k_2^2$ where $a = M\phi(s_i, a_i) - \gamma\phi(s_i^0, \pi)$, $b = g_\phi(s_i, a_i) - \gamma\phi(s_i^0, \pi)$. The envelope is $jX_ij = jX_ij = ka k_2^2 + kb k_2^2 \leq (1+\gamma)^2 + 8$. Further, observe that $ka k_2^2 - kb k_2^2 = ha + b$, $a = bi = ka + bk_2ka - bk_2$. So, the second moment is bounded $\mathbb{E}[X_i^2] \leq \mathbb{E}[ka + bk_2ka - bk_2]^2 \leq (1+3\gamma)^2 \mathbb{E}[kM\phi(s_i, a_i) - \gamma g_\phi(s_i, a_i)k_2^2] \leq 16\mathbb{E}[X_i]$, where we used Lemma E.3 to give us

$$\mathbb{E}[kM\phi(s_i, a_i) - \gamma g_\phi(s_i, a_i)k_2^2] = \mathbb{E}[kM\phi(s_i, a_i) - \gamma\phi(s_i^0, \pi)k_2^2 - kg_\phi(s_i, a_i) - \gamma\phi(s_i^0, \pi)k_2^2].$$

So, the error from the lemma is $\frac{24 \log(2/\delta)}{N}$. Now union bound over an ε -net of $\{M \in \mathbb{R}^{d \times d} : kMk_2 \leq 1\}$. Observe that

$$\begin{aligned}& \left| (kM\phi(s, a) - \gamma\phi(s^0, \pi)k_2^2 - kg_\phi(s, a) - \gamma\phi(s^0, \pi)k_2^2) - (k\tilde{M}\phi(s, a) - \gamma\phi(s^0, \pi)k_2^2 - kg_\phi(s, a) - \gamma\phi(s^0, \pi)k_2^2) \right| \\ & \left\| (M - \tilde{M})\phi(s_i, a_i) \right\|_2 \left\| M\phi(s_i, a_i) + \tilde{M}\phi(s_i, a_i) - 2\gamma\phi(s_i^0, \pi) \right\|_2 \\ & \left\| (M - \tilde{M}) \right\|_2 (2+2\gamma) \leq 4 \left\| (M - \tilde{M}) \right\|_2.\end{aligned}$$

Consider a $\frac{1}{N}$ -net (under kk_F) for $\mathcal{M} \in \mathbb{R}^{d \times d} : kMk_F \leq \frac{1}{d} \bar{d}g$, which requires $(1+2N\frac{1}{d}\bar{d})^{d^2}$ points since it is like the ℓ_2 for a d^2 -dimensional vector. This is a $\frac{1}{N}$ -net (under kk_2) for the subset $\{M \in \mathbb{R}^{d \times d} : kMk_2 \leq 1\}$ since $kMk_2 \leq kMk_F \leq \frac{1}{d} \bar{d}kMk_2$. The error from this ε -net approximation is at most $\frac{8}{N}$. Finally, union bound over \mathcal{F} . \square

Proof for ε_g . For a fixed $g \in G$, apply Lemma B.1 to $X_i = kg(s_i, a_i) - \gamma\phi^*(s_i^0, \pi)k_2^2 - kg_{\phi^*}(s, a) - \gamma\phi^*(s^0, \pi)k_2^2$, the excess regression loss. Under realizability Assumption 4.3, we have $\mathbb{E}[X_i] = kg - g_{\phi^*}k_2^2$, since

$$\begin{aligned}\mathbb{E}[X_i] & = \mathbb{E}\left[kg(s, a) - \gamma\phi^*(s^0, \pi)k_2^2 - kg_{\phi^*}(s, a) - \gamma\phi^*(s^0, \pi)k_2^2\right] \\ & = \mathbb{E}\left[kg(s, a) - g_{\phi^*}(s, a)k_2^2 + 2hg_{\phi^*}(s, a) - \gamma\phi^*(s^0, \pi), g(s, a) - g_{\phi^*}(s, a)\right]\end{aligned}$$

By definition of g_{ϕ^*} , we have $\mathbb{E}_{s^0 \sim P(s^0)} [g_{\phi^*}(s, a) - \gamma\phi^*(s^0, \pi)] = 0$, so,

$$= \mathbb{E}\left[kg(s, a) - g_{\phi^*}(s, a)k_2^2\right]$$

The envelope is $jX_ij \leq (2\gamma)^2$ and the second moment is bounded $\mathbb{E}[X_i^2] \leq \mathbb{E}\left[kg(s, a) + g_{\phi^*}(s, a) - 2\gamma\mathbb{E}_{a^0 \sim \pi(s^0)} [\phi^*(s^0, a^0)]k_2^2 kg(s, a) - g_{\phi^*}(s, a)k_2^2\right] \leq (4\gamma)^2 \mathbb{E}[X_i]$. So the error term from the lemma is $\frac{20\gamma^2 \log(2/\delta)}{N}$. Now union bound over G . \square

Infinite function classes When \mathcal{F} and G are infinite, we need to assume some metric entropy conditions. Then in the final step of the finite-class proofs above, we union bound on a well-chosen ε -net and collect an additional approximation error which is on the order of $O(1/N)$.

Assumption E.5. For $\mathcal{F} \supseteq \mathcal{F}, G \supseteq \mathcal{G}$, we assume there exists $p \geq 2$ such that $N(t, \mathcal{F}) \leq t^{-p}$, where the net is under the following distances,

$$\begin{aligned} d(\phi, \tilde{\phi}) &= E_D \left[\left\| \phi(s, a) - \tilde{\phi}(s, a) \right\|_2 \right] + E_\nu \left[\left\| \phi(s, a) - \tilde{\phi}(s, a) \right\|_2 \right] \\ &\quad + E_D \left[\left\| \phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi) \right\|_2 \right] + E_{s, a} \left[E_{D, s^\ell} \left[\left\| \phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi) \right\|_2 \right] \right] \\ &\quad + E_{s, a} \left[E_{\nu, s^\ell} \left[\left\| \phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi) \right\|_2 \right] \right] \\ d_G(g, \tilde{g}) &= \sqrt{E_D \left[k g(s_i, a_i) - \tilde{g}(s_i, a_i) k^2 \right]} + \sqrt{E_\nu \left[k g(s, a) - \tilde{g}(s, a) k^2 \right]} \end{aligned}$$

Note that this assumption is automatically satisfied for any $p > 0$ by VC classes (van der Vaart & Wellner, 1996, Theorem 2.6.4). Under this assumption, we have

$$\begin{aligned} \varepsilon_\rho &\leq \frac{d(1+W)^2(1+p) \log(WN/\delta)}{N} \\ \varepsilon_M &\leq \frac{d^2(1+p) \log(N/\delta)}{N} \\ \varepsilon_g &\leq \frac{\gamma^2 p \log(N/\delta)}{N}. \end{aligned}$$

Proof for ε_ρ . From before, we showed for a fixed $\phi \in \mathcal{F}$, we have $\varepsilon_\rho \leq \frac{d(1+W)^2 \log(WN/\delta)}{N}$. Observe that $\left| (\rho^\top \phi(s, a) - r(s, a))^2 - (\rho^\top \tilde{\phi}(s, a) - r(s, a))^2 \right| = \left| \rho^\top (\phi(s, a) - \tilde{\phi}(s, a)) (\rho^\top (\phi(s, a) + \tilde{\phi}(s, a)) + 2r(s, a)) \right| \leq 2W(1+W)k\phi(s, a) - \tilde{\phi}(s, a)k_2$, and so the difference of the loss with ϕ and the loss with $\tilde{\phi}$ is bounded by $2W(1+W)d(\phi, \tilde{\phi})$. Now union bound over a $\frac{1}{N}$ -net, which requires $O(N^p)$ points by Assumption E.5. The approximation error using the net is $\frac{2W(1+W)}{N}$. \square

Proof for ε_M . From before, we showed for a fixed $\phi \in \mathcal{F}$, we have $\varepsilon_M \leq \frac{d^2 \log(N/\delta)}{N}$. Observe that

$$\begin{aligned} &\left| (kM\phi(s, a) - \gamma\phi(s^\ell, \pi)k^2 - kg_\phi(s, a) - \gamma\phi(s^\ell, \pi)k^2) - (kM\tilde{\phi}(s, a) - \gamma\tilde{\phi}(s^\ell, \pi)k^2 - kg_{\tilde{\phi}}(s, a) - \gamma\tilde{\phi}(s^\ell, \pi)k^2) \right| \\ &\quad kM(\phi(s, a) - \tilde{\phi}(s, a)) - \gamma(\phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi))k \\ &\quad + kg_\phi(s, a) - g_{\tilde{\phi}}(s, a) - \gamma(\phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi))k \\ &\quad + kg_\phi(s, a) + g_{\tilde{\phi}}(s, a) - \gamma(\phi(s^\ell, \pi) + \tilde{\phi}(s^\ell, \pi))k \\ &\quad \left(k\phi(s, a) - \tilde{\phi}(s, a)k + \gamma k\phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi)k \right) (2 + 2\gamma) \\ &\quad + \left(kg_\phi(s, a) - g_{\tilde{\phi}}(s, a)k + \gamma k\phi(s^\ell, \pi) - \tilde{\phi}(s^\ell, \pi)k \right) (2 + 2\gamma) \end{aligned}$$

Using closed form solution for g_ϕ ,

$$16d(\phi, \tilde{\phi}).$$

Now union bound over a $\frac{1}{N}$ -net, which requires $O(N^p)$ points by Assumption E.5. The approximation error using the net is at most $\frac{16}{N}$. \square

Proof for ε_g . From before, we showed for a fixed $g \in G$, we have $\varepsilon_g \leq \frac{\gamma^2 \log(1/\delta)}{N}$. Observe that

$$\begin{aligned} & |kg(s, a) - \gamma\phi^*(s^\ell, \pi)k^2 - k\tilde{g}(s, a) - \gamma\phi^*(s^\ell, \pi)k^2| \\ & \quad |kg(s, a) - \tilde{g}(s, a)k| |kg(s, a) + \tilde{g}(s, a)| + 2\gamma\phi^*(s^\ell, \pi)k \\ & \quad (2 + 2\gamma)d_G(g, \tilde{g}). \end{aligned}$$

Now union bound over a $\frac{1}{N}$ -net, which requires $O(N^p)$ points by Assumption E.5. The approximation error using the net is at most $\frac{4}{N}$. \square

E.3. Main Results

Lemma E.6. *Suppose $\phi^* \in \mathcal{L}$ is Linear BC. Suppose Assumption 4.3 if transitions are stochastic. Moreover, suppose ϕ^* is feasible in the bilevel optimization (and so $\hat{L}_{lbc}(\hat{\phi}) = \hat{L}_{lbc}(\phi^*)$). Then, w.p. $1 - 5\delta$,*

$$L_{lbc}(\hat{\phi}) \leq \frac{24d(1+W)^2 \log(4Wj/N/\delta) + 128d^2 \log(4j/N/\delta) + 40\gamma^2 \log(2jGj/\delta)}{N}$$

Proof.

$$\begin{aligned} & L_{lbc}(\hat{\phi}) \\ & = \mathbb{E}_\nu \left[(\hat{\rho}_\phi | \hat{\phi}(s, a) - r(s, a))^2 \right] + \mathbb{E}_{\nu, P} \left[k\hat{M}_\phi | \hat{\phi}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] = \mathbb{E}_{\nu, P} \left[kg_{\hat{\phi}}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] \end{aligned}$$

Since $\hat{\rho}_\phi, \hat{M}_\phi$ are minimizers under ν ,

$$\mathbb{E}_\nu \left[(\hat{\rho}_\phi | \hat{\phi}(s, a) - r(s, a))^2 \right] + \mathbb{E}_{\nu, P} \left[k\hat{M}_\phi | \hat{\phi}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] = \mathbb{E}_{\nu, P} \left[kg_{\hat{\phi}}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right]$$

By the ρ, M parts of Lemma E.4,

$$2\mathbb{E}_D \left[(\hat{\rho}_\phi | \hat{\phi}(s, a) - r(s, a))^2 \right] + 2\mathbb{E}_D \left[k\hat{M}_\phi | \hat{\phi}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] = 2\mathbb{E}_D \left[kg_{\hat{\phi}}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] + 2\varepsilon_\rho + 2\varepsilon_M$$

Since \hat{g}_ϕ minimizes under D ,

$$2\mathbb{E}_D \left[(\hat{\rho}_\phi | \hat{\phi}(s, a) - r(s, a))^2 \right] + 2\mathbb{E}_D \left[k\hat{M}_\phi | \hat{\phi}(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] = 2\mathbb{E}_D \left[k\hat{g}_\phi(s, a) - \gamma\hat{\phi}(s^\ell, \pi)k_2^2 \right] + 2\varepsilon_\rho + 2\varepsilon_M$$

By optimality of $\hat{\phi}$ under \hat{L}_{lbc} ,

$$\begin{aligned} & 2\mathbb{E}_D \left[(\hat{\rho}_\phi | \phi^*(s, a) - r(s, a))^2 \right] + 2\mathbb{E}_D \left[k\hat{M}_\phi | \phi^*(s, a) - \gamma\phi^*(s^\ell, \pi)k_2^2 \right] = 2\mathbb{E}_D \left[k\hat{g}_{\phi^*}(s, a) - \gamma\phi^*(s^\ell, \pi)k_2^2 \right] \\ & + 2\varepsilon_\rho + 2\varepsilon_M \end{aligned}$$

By the G part of Lemma E.4, we have $\mathbb{E}_D \left[kg_{\phi^*}(s, a) - \phi^*(s^\ell, \pi)k_2^2 - k\hat{g}_{\phi^*}(s, a) - \phi^*(s^\ell, \pi)k_2^2 \right] \leq \frac{1}{2}\mathbb{E}_\nu \left[kg_{\phi^*} - \hat{g}_{\phi^*}k_2^2 \right] + \varepsilon_g = \varepsilon_g$. Then, using the optimality of $\hat{\rho}$ and \hat{M} under D ,

$$\begin{aligned} & 2\mathbb{E}_D \left[(\hat{\rho}_\phi | \phi^*(s, a) - r(s, a))^2 \right] + 2\mathbb{E}_D \left[k\hat{M}_\phi | \phi^*(s, a) - \gamma\phi^*(s^\ell, \pi)k_2^2 \right] = 2\mathbb{E}_D \left[kg_{\phi^*}(s, a) - \phi^*(s^\ell, \pi)k_2^2 \right] \\ & + 2\varepsilon_\rho + 2\varepsilon_M + 2\varepsilon_g \end{aligned}$$

By the ρ, M parts of Lemma E.4,

$$\begin{aligned} & 3\mathbb{E}_\nu \left[(\hat{\rho}_\phi | \phi^*(s, a) - r(s, a))^2 \right] + 3\mathbb{E}_{\nu, P} \left[k\hat{M}_\phi | \phi^*(s, a) - \gamma\phi^*(s^\ell, \pi)k_2^2 \right] = 3\mathbb{E}_\nu \left[kg_{\phi^*}(s, a) - \phi^*(s^\ell, \pi)k_2^2 \right] \\ & + 4\varepsilon_\rho + 4\varepsilon_M + 2\varepsilon_g \end{aligned}$$

By Assumption 4.3 and Lemma E.3,

$$= 3L_{bc}(\phi^*) + 4\varepsilon_\rho + 4\varepsilon_M + 2\varepsilon_g$$

By assumption that ϕ^* is Linear BC and Proposition 4.2,

$$= 4\varepsilon_\rho + 4\varepsilon_M + 2\varepsilon_g.$$

□

We now prove Theorem 4.4, in the general stochastic case. The deterministic transitions case is subsumed by ignoring the minimization over G , i.e. setting the complexity term of G to zero.

Theorem 4.4. *Assume Assumption 4.1 (and Assumption 4.3 if the system is stochastic). Let $C_2 := \frac{96 \log^{1/2}(j) + 4 \rho_{\bar{d}} + 4 \log^{1/2}(1/\delta)}{\beta/4}$. If $N \geq C_2^2$, then for any $\delta \geq (0, 1)$, w.p. at least $1 - \delta$, we have*

1. $\hat{\phi}$ satisfies $\hat{\varepsilon}$ -approximate Linear BC with

$$\hat{\varepsilon} = \frac{13d(1+W)^2 \log^{1/2}(4Wj/jN/\delta)}{\rho_{\bar{N}}} + \frac{7\gamma(1+W) \log^{1/2}(2jGj/\delta)}{\rho_{\bar{N}}},$$

2. $\lambda_{\min}(\hat{\phi}) \geq \beta/4$.

If transitions are deterministic, treat $\log(jGj) = 0$.

Proof of Theorem 4.4. First, by our assumption that $\frac{\rho_{\bar{N}}}{N} \geq 4(96\kappa(\cdot) + 4\rho_{\bar{d}} + 4 \log^{1/2}(1/\delta))/\beta$, Lemma E.2 implies that w.p. at least $1 - \delta$, we have $\sup_{\phi \geq 2} |\lambda_{\min}(\phi) - \lambda_{\min}(\hat{\phi})| \leq \beta/4$. Under this high probability event, we have two important consequences:

1. ϕ^* is feasible in Equation (6), since $\lambda_{\min}(\hat{\phi}) \geq \lambda_{\min}(\phi^*) - \beta/4 \geq \beta(1 - 1/4) \geq \beta/2$. In particular, this means $\hat{L}_{bc}(\hat{\phi}) \geq \hat{L}_{bc}(\phi^*)$.
2. The covariance of $\hat{\phi}$ has lower-bounded eigenvalues, since $\lambda_{\min}(\hat{\phi}) \geq \lambda_{\min}(\hat{\phi}) - \beta/4 \geq \beta(1/2 - 1/4) \geq \beta/4$.

Now, apply Lemma E.6 to bound $L_{bc}(\hat{\phi})$, so w.p. $1 - 5\delta$,

$$L_{bc}(\hat{\phi}) \leq \frac{24d(1+W)^2 \log(4Wj/jN/\delta) + 128d^2 \log(4j/jN/\delta) + 40\gamma^2 \log(2jGj/\delta)}{N}$$

By Lemma D.1, we have that $\hat{\phi}$ is $\hat{\varepsilon}$ -approximately Linear BC, with parameter

$$\hat{\varepsilon} = (1+W) \sqrt{\frac{24d(1+W)^2 \log(4Wj/jN/\delta) + 128d^2 \log(4j/jN/\delta) + 40\gamma^2 \log(2jGj/\delta)}{N}} + \frac{13d(1+W)^2 \log^{1/2}(4Wj/jN/\delta)}{\rho_{\bar{N}}} + \frac{7\gamma(1+W) \log^{1/2}(2jGj/\delta)}{\rho_{\bar{N}}}$$

Finally, we remark that the \widehat{W} for $\widehat{\phi}$ (in the approximately Linear BC case) is upper bounded by a polynomial in W^* in the assumed exact Linear BC of ϕ^* . Consider our assumption that ϕ^* is exactly Linear BC with $W^* = W$ (use \star to highlight that it is the W in the assumption, which we now show matches the W in the result). Then, by Proposition 4.2, $\mathcal{Q}M^*$ with $kM^*k_2 \sqrt{1 + \frac{k\rho^*k_2^2}{W^{\star 2}}}$. Hence, it suffices to minimize over this smaller ball for \widehat{M} , so that $k\widehat{M}k_2 \leq kM^*k_2$. Now, take the smallest possible W in Lemma D.1, so that

$$\begin{aligned} \widehat{W} &= \frac{k\rho^*k_2}{1 - kM^*k_2} \\ &= \frac{k\rho^*k_2}{1 - \sqrt{1 + \frac{k\rho^*k_2^2}{W^{\star 2}}}} \\ &= \frac{k\rho^*k_2 (1 + \sqrt{1 + \frac{k\rho^*k_2^2}{W^{\star 2}}})}{\frac{k\rho^*k_2^2}{W^{\star 2}}} \\ &= \frac{2W^{\star 2}}{k\rho^*k_2}, \end{aligned}$$

which is a polynomial in W^* . □

Our end-to-end result is deduced by chaining our LSPE theorem and the above theorem together.

Theorem 4.5. *Under Assumption 4.1 (and Assumption 4.3 if $P(s, a)$ is stochastic). Let $C_2, \widehat{\varepsilon}$ be as defined in Theorem 4.4. If $N \geq C_2^2$, then we have for any $\delta \geq (0, 1/2)$, w.p. at least $1 - 2\delta$, for all distributions p_0 ,*

$$\begin{aligned} \left| V_{p_0}^{\pi_e} - \mathbb{E}_{s \sim p_0} [\widehat{f}_K(s, \pi_e)] \right| &\leq \frac{\gamma^{K/2}}{1 - \gamma} + \frac{4 \sqrt{\left\| \frac{dd_{p_0}^{\pi_e}}{d\nu} \right\|_1}}{(1 - \gamma)^2} \widehat{\varepsilon} \\ &+ \frac{960\beta^{-1/2}(1 + W)d \log(N) \sqrt{\log(10/\delta)}}{(1 - \gamma)^2 N}. \end{aligned}$$

Proof of Theorem 4.5. We first apply Theorem 4.4 to see that $\widehat{\phi}$ satisfies the two properties needed for LSPE. It is indeed approximately Linear BC, with $\widehat{\varepsilon}$ specified in the theorem, and also has coverage (i.e. $\lambda_{\min}(\widehat{\phi}) \geq \beta/4$). Using these two facts, and on a separate independent dataset D_2 (needs to be a separate dataset since $\widehat{\phi}$ is data-dependent), we run LSPE and directly apply Theorem 3.3 for the result. □

F. Implementation Details

Here we detail all environment specifications and hyperparameters used in the main text.

F.1. Dataset Details

Using the publicly released implementation for DrQ-v2, we trained high quality target policies and saved checkpoints for offline behavior datasets. We refer the readers to Yarats et al. (2021) for exact hyperparameters.

F.2. Environment Details

Following the standards used by DrQ-v2 (Yarats et al., 2021), all environments have a maximum horizon length of 500 timesteps. This is achieved by the behavior/target policy having an action repeat of 2 frames. Furthermore, each state is 3 stacked frames that are each 84×84 dimensional RGB images (thus $9 \times 84 \times 84$).

Learning Bellman Complete Representations for Offline Policy Evaluation

Task	Target Performance	Behavior Performance
Finger Turn Hard	927	226 (24%)
Cheetah Run	758	192 (25%)
Quadruped Walk	873	236 (27%)
Humanoid Stand	827	277 (33%)

Table 1. Performance for target and behavior policies used to collect evaluation and offline datasets respectively.

Task	Action Space Dimension	Task Traits	Reward Type
Finger Turn Hard	2	turn	sparse
Cheetah Run	6	locomotion	dense
Quadruped Walk	12	locomotion	dense
Humanoid Stand	21	stand	dense

Table 2. Task descriptions, action space dimension, and reward type for each tested environment.

E.3. Representation Architecture and Hyperparameter Details

We adopt the same network architecture as DrQ-v2’s critic, first introduced in SAC-AE (Yarats et al., 2019). More specifically, to process pixel input, we have a 5 layer ConvNet with 3×3 kernels and 32 channels with ReLU activations. The first convolutional layer has a stride of 2 while the rest has stride 1. The output is fed through a single fully connected layer normalized by LayerNorm. Finally, there is a tanh nonlinearity on the outputted 50 dimensional state-representation. The action is then concatenated with this output and fed into a 4-layer MLP all with ReLU activations.

Hyperparameter	Value
Feature Dimension	512
Weight Initialization	orthogonal init.
Optimizer	Adam
Learning Rate	$1 \cdot 10^{-5}$
Batch Size	2048
Training Epochs	200
τ (target)	0.005
λ_{Design}	$5 \cdot 10^{-6}$

Table 3. Hyperparameters used for BCRL

F.4. Benchmarks and Metrics

Modifications to CURL: Originally, CURL only does contrastive learning between the image states with data augmentation. For OPE, apply the same CURL objective to the state-action feature detailed in the previous section. Note we also train CURL with the same random cropping image augmentations presented by the authors. Finally, since we are not interleaving the representation learning with SAC, we do not have a Q prediction head.

Modifications to SPR: We use the same image encoder as our features for SPR. The main difference is in the architecture of the projection layers where we implement as 3-layer mlp with ReLU activations. Note that these are additional parameters that neither CURL nor BCRL require. Finally, similarly to CURL, we do not have an additional Q-prediction head.

Spearman Ranking Correlation Metric: This rank correlation measures the correlation between the ordinal rankings of the value estimates and the ground truth returns. As defined in Fu et al. (2021), we have for policies $1, 2, \dots, N$, true returns $V_{1:N}$, and estimated returns $\hat{V}_{1:N}$:

$$\text{Ranking Correlation} = \frac{\text{Cov}(V_{1:N}, \hat{V}_{1:N})}{\sigma(V_{1:N}) \sigma(\hat{V}_{1:N})}$$