

---

# Primal-Dual Spectral Representation for Off-policy Evaluation

---

Yang Hu<sup>\*,1</sup>   Tianyi Chen<sup>\*,2</sup>   Na Li<sup>1</sup>   Kai Wang<sup>2</sup>   Bo Dai<sup>2</sup>  
<sup>1</sup>Harvard University   <sup>2</sup>Georgia Institute of Technology   <sup>\*</sup>Equal contribution.

## Abstract

Off-policy evaluation (OPE) is one of the most fundamental problems in reinforcement learning (RL) to estimate the expected long-term payoff of a given target policy with *only* experiences from another behavior policy that is potentially unknown. The distribution correction estimation (DICE) family of estimators have advanced the state of the art in OPE by breaking the *curse of horizon*. However, the major bottleneck of applying DICE estimators lies in the difficulty of solving the saddle-point optimization involved, especially with neural network implementations. In this paper, we tackle this challenge by establishing a *linear representation* of value function and stationary distribution correction ratio, *i.e.*, primal and dual variables in the DICE framework, using the spectral decomposition of the transition operator. Such primal-dual representation not only bypasses the non-convex non-concave optimization in vanilla DICE, therefore enabling an computational efficient algorithm, but also paves the way for more efficient utilization of historical data. We highlight that our algorithm, SPECTRALDICE, is the first to leverage the linear representation of primal-dual variables that is both computation and sample efficient, the performance of which is supported by a rigorous theoretical sample complexity guarantee and a thorough empirical evaluation on various benchmarks.

## 1 INTRODUCTION

The past decade has witnessed the ubiquitous success of reinforcement learning (RL) across various domains. Despite the original rationale that RL agents should learn a reward-maximizing policy from continuous in-

teractions with the environment, there also exist a wide range of applicational scenarios where *online* interaction with the environment may be expensive, inefficient, risky, unethical, and/or even infeasible, examples of which include robotics (Kalashnikov et al., 2018; Kahn et al., 2018), autonomous driving (Shi et al., 2021; Fang et al., 2022), healthcare (Jagannatha et al., 2018; Gottesman et al., 2018), education (Mandel et al., 2014; Slim et al., 2021), dialogue systems (Jaques et al., 2019; Jiang et al., 2021) and recommendation systems (Li et al., 2011; Chen et al., 2019). These application scenarios motivate the study of *offline* RL, where the learning agent only has access to historical data collected by a separate behavior policy.

Off-policy evaluation (OPE) is one of the most fundamental problems in offline RL that aims at estimating the expected cumulative reward of a given target policy using only historical data collected by a different, potentially unknown behavior policy. In the past decade, various off-policy performance estimators have been proposed (Hanna et al., 2019; Xie et al., 2019; Jiang and Li, 2016; Foster et al., 2021). However, these estimators generally suffer from the *curse of horizon* (Liu et al., 2018)—step-wise variances accumulate in a multiplicative way, resulting in prohibitively high trajectory variances and thus unreliable estimators. The recently proposed Distribution Correction Estimation (DICE) family of estimators have advanced the state of the art in OPE, leveraging the primal-dual formulation of policy evaluation for a saddle-point optimization approach that directly estimates the stationary distribution correction ratio, and hence breaking the curse of horizon (Nachum et al., 2019a,b).

Nevertheless, as systems scale up in terms of the size of state-action spaces, the saddle-point optimization in the formulation of DICE estimators become increasingly challenging to solve. Such *curse of dimensionality* is common for RL methods in general, and people have been working to alleviate the computational burden by exploiting function approximators. However, many known function approximators require additional assumptions to ensure computational and statistical properties (Gordon, 1995; Jiang et al., 2017; Chen and Jiang,

2019; Zhan et al., 2022; Katdare et al., 2023; Che et al., 2024), which may not be easily satisfiable in practice. Moreover, the induced optimization upon function approximators may be difficult to solve (Boyan and Moore, 1994; Baird, 1995; Tsitsiklis and Van Roy, 1996). In particular, under a generic neural network parametrization, computing the DICE estimator (Nachum and Dai, 2020) requires solving non-convex non-concave saddle-point optimizations, which is known to be NP-hard in theory and also yields unstable performance in practice, and is therefore regarded as intractable.

This dilemma brings up a very natural question:

*Can we design an OPE algorithm that is both **efficient** and **practical**?*

By “efficient” we mean its statistical complexity avoids an exponential dependence on both the length of history and the dimension of state-action spaces, *i.e.*, eliminating both *curse of horizon* and *curse of dimensionality*. By “practical” we mean the algorithm is free from unstable saddle-point optimizations and can be easily implemented and applied in practical settings.

In this paper, we provide an *affirmative* answer to this question by revealing a novel linear structure encapsulating both  $Q$ -functions and distribution correction ratios via a spectral representation of the transition operator, which has many nice properties to enable efficient representation learning and off-policy evaluation.

**Contributions.** Specifically, the contributions of this paper can be summarized as follows:

- We propose a novel *primal-dual spectral representation* of the state-action transition operator, which makes both the  $Q$ -function and the stationary distribution correction ratio (*i.e.*, the primal and dual variables in DICE) linearly representable in the primal/dual feature spaces, and thus enhances the tractability of the corresponding DICE estimator.
- We design SPECTRALDICE, an off-policy evaluation algorithm based on our primal-dual spectral representation, which bypasses the non-convex non-concave saddle-point optimization in vanilla DICE with generic neural network parameterization, and also makes efficient use of historical data. As far as we are concerned, our algorithm is the first to leverage the linear representation of both primal and dual variables that is *computation and sample efficient*.
- The performance of the SPECTRALDICE algorithm is justified both theoretically with a rigorous sample complexity guarantee and empirically by a thorough evaluation on various RL benchmarks.

## 1.1 Related Work

**Off-Policy Evaluation (OPE).** Off-policy evaluation has long been an active field of RL research. In

the case where the behavior policy is known, various off-policy performance estimators have been proposed, including direct method (DM) estimators (Antos et al., 2008; Le et al., 2019), importance sampling (IS) estimators (Hanna et al., 2019; Xie et al., 2019), doubly-robust (DR) estimators (Dudík et al., 2011; Jiang and Li, 2016; Foster et al., 2021) and other mixed-type estimators (Thomas and Brunskill, 2016; Kallus and Uehara, 2020; Katdare et al., 2023), which generally suffer from the *curse of dimension*. In an effort to settle this issue, there is also abundant literature on estimating the correction ratio of the stationary distribution (Liu et al., 2018; Uehara et al., 2020), among which the distribution correction estimation (DICE) family of estimators are the state of the art that leverage a novel primal-dual formulation of OPE to eliminate the curse of horizon, and in the meantime, allow unknown behavior policies (Nachum et al., 2019a,b; Yang et al., 2022; Zhang et al., 2020; Nachum and Dai, 2020). However, as discussed above, the induced saddle-point optimization becomes unstable with neural networks, impeding the practical application of DICE estimators.

**Spectral Representation in MDPs.** Spectral decomposition of the transition kernel is known to induce a linear structure of  $Q$ -functions, which enables the design of provably efficient algorithms assuming known (primal) spectral feature maps (Jin et al., 2020; Yang and Wang, 2020; Ren et al., 2022b). These algorithms break the curse of dimensionality in the sense that their computation or sample complexity is independent of the size of the state-action space, but rather, only depends polynomially on the feature space dimension, the intrinsic dimension of the problem.

With the growing interest in spectral structures of MDPs, representation learning for RL has recently attracted much theory-oriented attention in the online setting (Agarwal et al., 2020; Uehara et al., 2021). Practical representation-based online RL algorithms have been designed via kernel techniques (Ren et al., 2022c, 2023), latent variable models (Ren et al., 2022a; Zhang et al., 2023), contrastive learning (Qiu et al., 2022; Zhang et al., 2022a), and diffusion score matching (Shribak et al., 2024). Recently, a unified representation learning framework is proposed from a novel viewpoint that leverages the spectral decomposition of the transition operator (Ren et al., 2022b).

Spectral representations have also been exploited in the offline setting (Uehara and Sun, 2021; Ni et al., 2021; Chang et al., 2022), where the temporal difference algorithm is applied in the linear space induced by the primal spectral feature for estimating  $Q$ -functions. The linear structure of the occupancy measure induced by the dual spectral feature is recently utilized in Huang et al. (2023), which leads to an offline RL algorithm

for stationary density ratio estimation. Although the algorithm is theoretically sound, the stationary density ratio breaks the linearity in occupancy, and hence the algorithm is not computationally efficient. As far as we know, there is no such offline RL algorithm that efficiently utilizes both primal and dual representations.

## 2 PRELIMINARIES

**Notations.** Denote by  $\|\cdot\|_p$  the  $p$ -norm of vectors or the  $L^p$ -norm of functionals, and by  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$  the Euclidean inner product of vectors  $\mathbf{x}$  and  $\mathbf{y}$ . Denote by  $\mathbb{E}_{d^{\mathcal{D}}}[\cdot]$  the empirically approximated expectation using samples from dataset  $\mathcal{D} \sim d^{\mathcal{P}}$ . Denote by  $\Delta(S)$  the set of distributions over set  $S$ , the element of which shall be regarded as densities whenever feasible. Denote the indicator function by  $\mathbb{I}\{\cdot\}$ . Write  $[n] := \{1, \dots, n\}$  for  $n \in \mathbb{Z}_+$ . Regard  $f(n) \lesssim g(n)$  as  $f(n) = O(g(n))$ .

**Markov Decision Processes (MDPs).** We consider an *infinite-horizon* discounted Markov decision process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \mu_0, \gamma)$ , where  $\mathcal{S}$  is the (possibly infinite) state space,  $\mathcal{A}$  is the (possibly infinite) action space;  $\mathbb{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition kernel,  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function;  $\mu_0 \in \Delta(\mathcal{S})$  is the initial state distribution, and  $\gamma \in (0, 1)$  is the reward discount factor, so that the discounted cumulative reward can be defined as  $\sum_{t=0}^{\infty} \gamma^t r_t$ . We consider *stationary Markovian policies*  $\Pi := \{\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})\}$  that admit an action distribution depending on the current state only. Given any policy  $\pi \in \Pi$ , let  $\mathbb{E}_{\pi, \mathbb{P}}[\cdot]$  denote the expectation over the trajectory governed by  $\pi$  and  $\mathbb{P}$  (possibly under prescribed initial conditions). Let  $d_{\mathbb{P}}^{\pi}(\cdot, \cdot) \in \Delta(\mathcal{S} \times \mathcal{A})$  denote the (*stationary*) *state-action occupancy measure* under policy  $\pi$ , i.e., the normalized discounted probability of visiting  $(s, a)$  in a trajectory induced by policy  $\pi$ , defined by

$$d_{\mathbb{P}}^{\pi}(s, a) = (1 - \gamma) \mathbb{E}_{\pi, \mathbb{P}} \left[ \sum_{t=0}^{\infty} \gamma^t \mathbb{I}\{s_t = s, a_t = a\} \right].$$

Similarly, let  $d_{\mathbb{P}}^{\pi}(\cdot) \in \Delta(\mathcal{S})$  denote the *state occupancy measure* subject to the relation  $d_{\mathbb{P}}^{\pi}(s, a) = d_{\mathbb{P}}^{\pi}(s) \pi(a|s)$ . Further, define the state/state-action value functions (a.k.a.  $V$ - and  $Q$ -functions) as follows:

$$V_{\mathbb{P}}^{\pi}(s) := \mathbb{E}_{\pi, \mathbb{P}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right],$$

$$Q_{\mathbb{P}}^{\pi}(s, a) := \mathbb{E}_{\pi, \mathbb{P}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

In this way, the value of policy  $\pi$  in  $\mathcal{M}$  is defined by

$$\begin{aligned} \rho_{\mathbb{P}}(\pi) &:= (1 - \gamma) \mathbb{E}_{s \sim \mu_0} [V_{\mathbb{P}}^{\pi}(s)] \\ &= (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q_{\mathbb{P}}^{\pi}(s, a)], \end{aligned} \quad (1)$$

where the factor  $(1 - \gamma)$  is introduced for normalization. We omit the subscript  $\mathbb{P}$  when the context is clear.

*Remark 1.* In order to better illustrate how the proposed method works in MDPs with continuous state-action spaces, we abuse the notation a bit to regard  $\mathbb{P}$ ,  $\pi$  and  $d^{\pi}$  as *densities*. Parallel results for the discrete case can be analogously derived without difficulties.

**The Primal-Dual Characterization of  $\rho(\pi)$ .** Distribution Correction Estimation (DICE) (Nachum and Dai, 2020) is a primal-dual-based method that evaluates the value of a given target policy  $\pi$  in the offline setting, using the linear programming (LP) formulation of policy values (Puterman, 2014). Specifically, it is known that we can equivalently characterize  $\rho(\pi)$  defined in (1) by the following *primal LP*:

$$\begin{aligned} \min_{Q(\cdot, \cdot)} \quad & (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)], \\ \text{s.t.} \quad & Q(s, a) \geq r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')], \\ & \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \end{aligned} \quad (2)$$

Further, it can be shown that strong duality holds in (2), with Lagrangian multipliers exactly the state-action occupancy measures  $d^{\pi}(\cdot, \cdot)$ . We can therefore characterize  $\rho(\pi)$  by the following *primal-dual LP*:

$$\begin{aligned} \min_{Q(\cdot, \cdot)} \max_{d(\cdot, \cdot)} \quad & (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \\ & \mathbb{E}_{(s, a) \sim d^{\pi}(\cdot, \cdot)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right]. \end{aligned} \quad (3)$$

We highlight that this primal-dual LP formulation is favored in the offline RL setting in that historical experiences can be utilized to empirically approximate the expectations in (3) after some simple change-of-variables. In particular, for any measurable function  $f(s, a)$ , the importance sampling (IS) estimator for the expected value of  $f(s, a)$  is given by

$$\mathbb{E}_{(s, a) \sim d^{\pi}} [f(s, a)] = \mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[ \frac{d^{\pi}(s, a)}{d^{\pi_b}(s, a)} \cdot f(s, a) \right], \quad (4)$$

where  $\zeta(s, a) := \frac{d^{\pi}(s, a)}{d^{\pi_b}(s, a)}$  is known as the *stationary distribution correction ratio* for dataset  $\mathcal{D} \sim d^{\mathcal{D}}$ .

The DICE family estimators (Nachum et al., 2019a; Zhang et al., 2022b; Dai et al., 2020) is designed by plugging the IS expectation estimator (4) into (3), such that the stationary distribution correction ratio  $\zeta(\cdot, \cdot)$  is parameterized along with the  $Q$ -function to formulate an optimization, with various regularization available (Yang et al., 2020). It is evident that the DICE family estimators are applicable to the offline RL setting with unknown behavior policy.

**Spectral Representation.** We can always perform spectral decomposition of the dynamic operator to obtain a spectral representation of *any* MDP (Ren et al., 2022b). In particular, *low-rank MDPs* refer to such MDPs with intrinsic finite-rank spectral representation structures that enable scalable RL algorithms, and are thus of theoretical interest (Yao et al., 2014; Jin et al., 2020). Formally,  $\mathcal{M}$  is said to be a *low-rank MDP* if there exists a *primal* feature map  $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$  and *dual* features  $\tilde{\mu} : \mathcal{S} \rightarrow \mathbb{R}^d$ ,  $\theta_r \in \mathbb{R}^d$ , such that  $\mathbb{P}(s'|s, a) = \langle \phi(s, a), \tilde{\mu}(s') \rangle$ ,  $r(s, a) = \langle \phi(s, a), \theta_r \rangle$ , for any  $s, s' \in \mathcal{S}$ ,  $a \in \mathcal{A}$ . Here both the primal feature  $\phi$  and the dual features  $\tilde{\mu}, \theta_r$  are assumed to be unknown,

and thus must be learned from data (Agarwal et al., 2020; Uehara et al., 2021).

We would like to highlight that the assumption of factorizability should not be deemed as restrictive here. In fact, every transition kernel is technically “factorizable” in a sufficiently large infinite-dimensional space, though the factorization is not useful without a sufficiently low intrinsic dimension  $d$ . It is known in literature that, under very mild assumptions, approximate representations can be obtained by truncation of infinite-dimensional representations (Ren et al., 2023) or by latent variable models (Ren et al., 2022a) with bounded approximation error. For technical simplicity, we directly assume existence of such low-rank structures here, which is also standard in linear/low-rank MDP literature (Jin et al., 2020; Yang and Wang, 2020; Agarwal et al., 2020; Uehara et al., 2021).

Unfortunately, it is revealed in Zhang et al. (2022a); Ren et al. (2022b) that learning the features of a low-rank MDP is difficult from the unnormalized density fitting point of view. To settle this tractability issue, the above papers propose a reparameterization of the dual feature as  $\tilde{\mu}(\cdot) = q(\cdot)\mu(\cdot)$ , where we introduce an auxiliary distribution  $q(\cdot) \in \Delta(\mathcal{S})$  that will be specified later. Therefore, we will stick to the following spectral decomposition of the transition kernel in this paper:

$$\mathbb{P}(s'|s, a) = \langle \phi(s, a), q(s')\mu(s') \rangle, \quad \forall s, a, s'. \quad (5)$$

Under such reparameterization, it has been shown that the spectral representation can be learned efficiently.

Additionally, we also assume the initial distribution  $\mu_0$  to be linearly representable in the dual feature space.

**Assumption 1** (initial representation). There exists  $\omega_0 \in \mathbb{R}^d$ , such that  $\mu_0(s) = q(s)\langle \mu(s), \omega_0 \rangle$ ,  $\forall s \in \mathcal{S}$ .

**Off-Policy Evaluation (OPE).** We consider a setting where we are given  $\mathcal{D} = \{(s_i, a_i, s'_i) \mid i \in [N]\}$ , an offline dataset of  $N$  historical transitions, sampled by certain *behavior policy*  $\pi_b$  that could be unknown. The objective is to estimate the expected cumulative rewards  $\rho(\pi)$  of a different *target policy*  $\pi$ .

For satisfactory performance, it is important that the behavior policy provides sufficient data coverage for the frequent transitions experienced by policy  $\pi$ . Specifically, we assume the occupancy ratio between  $\pi$  and  $\pi_b$  satisfies the following regularity assumption.

**Assumption 2** (concentratability).  $\frac{d^\pi(s, a)}{d^{\pi_b}(s, a)} \leq C_\infty^\pi$ .

We point out that the concentratability assumption is standard in offline RL literature (Munos and Szepesvári, 2008; Chen and Jiang, 2019), and is also implicitly enforced in recent work like Huang et al. (2023) (see Definition 1 therein). We are aware that the coefficient  $C_\infty^\pi$  can potentially be translated into different feature-

related constants (Uehara et al., 2021), which does not change the asymptotics of sample complexity, yet only adds to the technical complexity. For clarity, we will stick to the simple Assumption 2 in this paper.

### 3 SPECTRALDICE: OPE USING PRIMAL-DUAL SPECTRAL REPRESENTATION

In this section, we first introduce a novel linear representation for the stationary distribution correction ratio using the *dual* spectral feature of transition kernel. We highlight that this linear structure, together with the known linear representation of  $Q$ -functions, helps to bypass the non-convex non-concave optimization required in the computation of DICE estimators, and also enables efficient utilization of historical data sampled by unknown behavior policies. Based on the above ideas, we present SPECTRALDICE, the proposed off-policy evaluation (OPE) algorithm using our primal-dual spectral representation.

#### 3.1 Primal-Dual Spectral Representation

We start by specifying the primal-dual spectral representation used in SPECTRALDICE. At first glance, it may seem natural to directly learn the spectral representation of  $\mathbb{P}$  as defined in (5). However, it turns out that this naive approach includes the target policy  $\pi$  in the linear representation of  $d^\pi(\cdot, \cdot)$ , which in turn induces a complicated representation for the stationary distribution correction ratio  $\zeta(\cdot, \cdot)$  (Huang et al., 2023), and thus, leads to an intractable optimization (3) for the computation of the DICE estimator.

The above challenge inspires us to properly reparameterize the spectral decomposition (5). Specifically, since we only work with a fixed target policy  $\pi$  for off-policy evaluation, we shall consider the following alternative representation of the state-action transition kernel  $\mathbb{P}^\pi(s', a'|s, a) := \mathbb{P}(s'|s, a)\pi(a'|s')$ :

$$\mathbb{P}^\pi(s', a'|s, a) = \left\langle \phi(s, a), q(s')\pi_b(a'|s') \underbrace{\frac{\pi(a'|s')}{\pi_b(a'|s')}}_{\mu^\pi(s', a')} \mu(s') \right\rangle. \quad (6)$$

Note that Assumption 2 guarantees a non-zero denominator when the numerator is non-zero. We refer to (6) as the *primal-dual spectral representation* of the state-action transition kernel  $\mathbb{P}^\pi$ , where  $\phi(\cdot, \cdot)$  and  $\mu^\pi(\cdot, \cdot)$  are still called *primal* and *dual* spectral features, respectively. The superscript  $\pi$  of the dual spectral feature emphasizes its dependence on the target policy.

The primal-dual spectral representation has several nice properties. In particular, we can show that the  $Q$ -function  $Q^\pi(s, a)$ , the state-action occupancy measure  $d^\pi(s, a)$ , and the stationary distribution correction ratio  $\zeta(s, a)$  can all be represented in linear forms using the



primal/dual features, as summarized below.

**Lemma 1.** *With primal-dual spectral representation (6), the  $Q$ -function  $Q^\pi(\cdot, \cdot)$  is linearly representable in the primal feature space with cofactor  $\theta_Q^\pi \in \mathbb{R}^d$ :*

$$Q^\pi(s, a) = \langle \phi(s, a), \theta_Q^\pi \rangle, \quad \forall s \in S, a \in \mathcal{A}. \quad (7)$$

*Further, under Assumption 1, the state-action occupancy measure  $d^\pi(\cdot, \cdot)$  is also linearly representable in the dual feature space with cofactor  $\omega_d^\pi \in \mathbb{R}^d$ :*

*$d^\pi(s, a) = q(s)\pi_b(a|s)\langle \mu^\pi(s, a), \omega_d^\pi \rangle$ ,  $\forall s \in S, a \in \mathcal{A}$ . Specifically, when the auxiliary distribution  $q(\cdot)$  is selected as the state-occupancy measure  $d^{\pi_b}(\cdot)$  of the behavior policy  $\pi_b$ , the stationary distribution correction ratio can also be linearly represented as:*

$$\zeta(s, a) = \frac{d^\pi(s, a)}{q(s)\pi_b(a|s)} = \langle \mu^\pi(s, a), \omega_d^\pi \rangle. \quad (8)$$

*Proof.* Note that the original dual feature in (5) can be restored by  $\mu(s') = \frac{\pi_b(a'|s')}{\pi(a'|s')} \mu^\pi(s', a')$  for any  $a' \in \mathcal{A}$ . Then by Bellman recursive equation we have:

$$\begin{aligned} Q^\pi(s, a) &= \langle \phi(s, a), \theta_r \rangle + \gamma \int V^\pi(s') \langle \phi(s, a), q(s') \mu(s') \rangle ds' \\ &= \left\langle \phi(s, a), \underbrace{\theta_r + \gamma \int V^\pi(s') q(s') \mu(s') ds'}_{\theta_Q^\pi} \right\rangle. \end{aligned}$$

Similarly, by the recursive property of  $d^\pi$  we have:

$$\begin{aligned} d^\pi(s, a) &= (1 - \gamma) \mu_0(s) \pi(a|s) + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \mathbb{P}^\pi(s, a | \tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a} \\ &= (1 - \gamma) q(s) \langle \pi_b(a|s) \mu^\pi(s, a), \omega_0 \rangle + \\ &\quad \gamma \left\langle q(s) \pi_b(a|s) \mu^\pi(s, a), \int d^\pi(\tilde{s}, \tilde{a}) \phi(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a} \right\rangle \\ &= \left\langle q(s) \pi_b(a|s) \mu^\pi(s, a), \right. \\ &\quad \left. \underbrace{(1 - \gamma) \omega_0 + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \phi(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a}}_{\omega_d^\pi} \right\rangle, \end{aligned}$$

where we use the initial representation (Assumption 1) and the fact that  $\pi(a|s) \mu(s) = \pi_b(a|s) \mu^\pi(s, a)$ . The representation of  $\zeta(\cdot, \cdot)$  is hence a direct corollary since  $q(s) \pi_b(a|s) = d^{\pi_b}(s, a)$  when  $q(\cdot) = d^{\pi_b}(\cdot)$ .  $\square$

Then, using the linear spectral representations of  $Q$  and  $\zeta$  in (7) and (8), we shall equivalently formulate the DICE estimator as follows.

**Corollary 2.** *With primal-dual spectral representation (6) where  $q(\cdot) \equiv d^{\pi_b}(\cdot)$ , under Assumption 1, we have*

$$\begin{aligned} \rho_\mathbb{P}(\pi) &= \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [\phi(s, a)^\top \theta_Q] \right. \\ &\quad \left. + \mathbb{E}_{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[ (\mu^\pi(s, a)^\top \omega_d) \cdot \right. \right. \\ &\quad \left. \left. (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q) \right] \right\}. \end{aligned} \quad (10)$$

The proof of Corollary 2 is deferred to Appendix B.1 due to limited space. We highlight that our new DICE

formulation (10) bears several benefits:

- **Offline data compatible.** The estimator is favorable for OPE since the expectation over the  $(s, a, s')$  transition pair can be effectively approximated by samples from the offline dataset  $\mathcal{D}$ , as long as the auxiliary distribution  $q(\cdot)$  is selected as the state occupancy measure  $d^{\pi_b}$  of the behavior policy  $\pi_b$  such that  $\Pr[(s, a, s') \in \mathcal{D}] = q(s)\pi_b(a|s)\mathbb{P}(s'|s, a)$ .
- **Optimization tractable.** Given (learned)  $\phi(s, a)$  and  $\mu^\pi(s, a)$ , the saddle-point optimization in (10) is convex-concave with respect to both  $\theta_Q$  and  $\omega_d$ , which perfectly bypasses the optimization difficulty in vanilla DICE estimators with neural-network-parameterized  $Q^\pi(\cdot, \cdot)$  and  $\zeta(\cdot, \cdot)$ . Meanwhile, compared to the counterpart obtained by directly applying the naive spectral representation (5) (details of which can be found in Appendix B.2), the proposed estimator (10) is tractable in that it is free of the policy ratio  $\frac{\pi(a|s)}{\pi_b(a|s)}$  that is unknown.

From now on, we will always regard  $q(\cdot) \equiv d^{\pi_b}(\cdot)$  for the aforementioned nice properties to hold.

### 3.2 Spectral Representation Learning

In the last section, we have elaborated on how to perform OPE using off-policy data given a primal-dual spectral representation. Now it only suffices to specify how to learn such a representation, which we regard as an abstract subroutine  $(\hat{\phi}, \hat{\mu}^\pi) \leftarrow \text{REPLEARN}(\mathcal{F}, \mathcal{D}, \pi)$ . Here  $\mathcal{F}$  denotes the collection of candidate representations. We highlight that our algorithm works with any representation learning method that has a bounded learning error, without any further requirements on the learning mechanism. Given a range of spectral representation learning methods available in literature (Zhang et al., 2022a; Ren et al., 2022b,a; Shribak et al., 2024), for the sake of clarity we only consider a few candidates here, while other methods may also be applicable:

1. **Ordinary Least Squares (OLS).** Inspired by Ren et al. (2022b), an OLS objective can be constructed as follows. Denote by  $\mathbb{Q}^\pi(s', a', s, a) := d^{\pi_b}(s, a) \mathbb{P}^\pi(s', a' | s, a)$  the joint distribution of state-action transitions under behavior policy  $\pi_b$ , based on which we plug in (6) to obtain  $\frac{\mathbb{Q}^\pi(s', a', s, a)}{\sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')}} = \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \phi(s, a)^\top \mu^\pi(s', a')$ , which further induces the following OLS objective:

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \int \left( \frac{\mathbb{Q}^\pi(s', a', s, a)}{\sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')}} - \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right)^2 ds da ds' da'$$

Therefore,  $(\hat{\phi}, \hat{\mu}^\pi)$  can be learned by solving (Ren et al., 2022b; HaoChen et al., 2021):

---

**Algorithm 1 SPECTRALDICE: Distribution Correction Estimation with Spectral Representation**


---

**Require:** Target policy  $\pi$ , off-policy dataset  $\mathcal{D}$ , function family  $\mathcal{F}$ .

- 1: Learn a spectral representation  $(\hat{\phi}, \hat{\mu}^\pi) \leftarrow \text{REPLEARN}(\mathcal{F}, \mathcal{D}, \pi)$ .
- 2: Plug in the spectral representation  $(\hat{\phi}, \hat{\mu}^\pi)$  to compute the following DICE estimator:

$$\hat{\rho}(\pi) = \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} \left[ \hat{\phi}(s, a)^\top \theta_Q \right] \right. \\ \left. + \mathbb{E}_{\substack{(s, a, s') \sim \mathcal{D}, \\ a' \sim \pi(\cdot|s')}} \left[ (\hat{\mu}^\pi(s, a)^\top \omega_d) (r(s, a) + \gamma \hat{\phi}(s', a')^\top \theta_Q - \hat{\phi}(s, a)^\top \theta_Q) \right] \right\}. \quad (9)$$

- 3: **return**  $\hat{\rho}(\pi)$
- 

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \left\{ \mathbb{E}_{(s, a) \sim d^{\pi_b}, (\tilde{s}', \tilde{a}') \sim d^{\pi_b}} \left[ (\hat{\phi}(s, a)^\top \hat{\mu}^\pi(\tilde{s}', \tilde{a}'))^2 \right] \right. \\ \left. - 2 \mathbb{E}_{(s, a) \sim d^{\pi_b}, (s', a') \sim \mathbb{P}^\pi(\cdot, \cdot|s, a)} \left[ \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right] \right\},$$

where the last term becomes a constant after expansion and is thus omitted. For practical implementation, we can use stochastic gradient descent to solve the above stochastic optimization problem.

2. **Noise-Contrastive Estimation (NCE).** NCE is a widely used method for contrastive representation learning in RL (Zhang et al., 2022a; Qiu et al., 2022). To learn  $(\hat{\phi}, \hat{\mu}^\pi)$ , we consider a binary contrastive learning objective (Qiu et al., 2022):

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \hat{E}_{(s, a) \sim d^{\pi_b}} \left[ \mathbb{E}_{(s', a') \sim \mathbb{P}^\pi(\cdot, \cdot|s, a)} \left[ \log \left( 1 + \frac{1}{\hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a')} \right) \right] \right] \\ + \mathbb{E}_{(s', a') \sim P_{\text{neg}}} \left[ \log \left( 1 + \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right) \right],$$

where  $P_{\text{neg}}$  is a negative sampling distribution.

Details of these representation learning methods along with their learning errors can be found in Appendix C.

### 3.3 SPECTRALDICE

With the two key components specified above, now we are ready to state SPECTRALDICE, the proposed offline policy evaluation (OPE) algorithm using spectral representations, as displayed in Algorithm 1.

Specifically, given a policy  $\pi$ , assuming access to an offline dataset  $(s, a, s') \sim \mathcal{D}$  sampled by the behavior policy  $\pi_b$ , we follow a two-step algorithm to evaluate the target policy  $\pi$  in an off-policy manner:

1. **Representation learning.** We may choose any representation learning method that comes with a bounded learning error as the REPLEARN subroutine, and the overall sample complexity will depend on this choice (see Section 4).
2. **DICE-based policy evaluation.** With the learned representation  $(\hat{\phi}, \hat{\mu}^\pi)$ , we use the primal-dual DICE estimator (9) to estimate the value of the target policy  $\pi$ . Note that the data distribution  $d^{\mathcal{D}}(s, a, s') = d^{\pi_b}(s) \pi_b(a|s) \mathbb{P}(s'|s, a)$  is exactly compatible with the formulation in (10).

*Remark 2* (Numerical considerations). It is known that

directly solving (9) leads to potential numerical instability issues due to the objective’s linearity in  $\theta_Q$  and  $\omega_d$  (Nachum et al., 2019b). Fortunately, it is shown in Yang et al. (2020) that certain regularization leads to strictly concave inner maximization while keeping the optimal *solution* unbiased (see Appendix B.3 for details). In our implementation, we append the following regularizer to the objective in (9):

$$-\lambda \mathbb{E}_{(s, a) \sim \mathcal{D}} [f(\hat{\mu}^\pi(s, a)^\top \omega_d)],$$

where  $f$  is a differentiable function with closed and convex Fenchel conjugate  $f_*$  (see Appendix E.1), and  $\lambda$  is a tunable constant. Furthermore, we also restrict  $\theta_Q$  and  $\omega_d$  in regions  $\Theta(\hat{\phi}) = \{\theta_Q \mid 0 \leq \hat{\phi}(s, a)^\top \theta_Q \leq \frac{1}{1-\gamma}\}$  and  $\Omega(\hat{\mu}^\pi) = \{\omega_d \mid \hat{\mu}^\pi(s, a)^\top \omega_d \leq C_\infty^\pi\}$ , respectively.

## 4 THEORETICAL GUARANTEE

In this section, we provide a rigorously theoretical analysis regarding the sample complexity of the proposed SPECTRALDICE algorithm. For the sake of technical conciseness, we make the following assumption on the candidate family  $\mathcal{F}$ . We argue that this is not a restrictive assumption, but rather, only helps to highlight the key contributions with simplified analysis.

**Assumption 3** (realizability). Assume a finite family  $\mathcal{F}$ , such that  $\langle \hat{\phi}(s, a), d^{\pi_b}(s', a') \hat{\mu}^\pi(s', a') \rangle$  is a valid state-action transition kernel for any  $(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}$ , and the ground-truth representation  $(\phi^*, \mu^{\pi, *}) \in \mathcal{F}$ .

*Remark 3* (Hypothesis class.). Here we assume  $\mathcal{F}$  to be finite for technical conciseness, which is consistent with literature (Agarwal et al., 2020; Uehara et al., 2021) to avoid analyzing the intrinsic complexity of hypothesis classes. It is expected that standard machine learning theory techniques can be used to adapt to infinite hypothesis classes with low intrinsic complexity. In addition, the realizability assumption  $(\phi^*, \mu^{\pi, *}) \in \mathcal{F}$  can be replaced by a (less restrictive) upper bound on realization error, such that the representation learning error  $\xi(|\mathcal{F}|, N, \delta)$  below will also include this error linearly. Since the overall evaluation error bound depends linearly on  $\xi(|\mathcal{F}|, N, \delta/2)$  (see Theorem 4), the additional realization error does not propagate over the steps, and thus will not deteriorate the bound.

**Representation Learning Error.** The key to subsequent analyses is to first bound the error of representation learning, which is of some theoretical interest by itself. Generally speaking, we expect *probably approximately correct* (PAC) bounds for representation learning in the following format.

**Claim 3.** *With probability at least  $1 - \delta$ , we have*

$$\mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi}} \left[ \left\| \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a) - \mathbb{P}^{\pi}(\cdot, \cdot | s, a) \right\|_1 \right] \leq \xi(|\mathcal{F}|, N, \delta),$$

where  $\hat{\mathbb{P}}^{\pi}(s', a' | s, a) := d_{\mathbb{P}}^{\pi}(s', a') \hat{\phi}(s, a)^{\top} \hat{\mu}^{\pi}(s')$ ,  $N$  is the number of samples in  $\mathcal{D}$ , and the upper bound  $\xi$  only depends on  $|\mathcal{F}|$ ,  $N$  and  $\delta$ .

We point out that, under certain regularity assumptions, the above claim can be proven for many spectral representation learning algorithms. Specifically, when REPLEARN is implemented by OLS or NCE, we can show that  $\xi(|\mathcal{F}|, N, \delta) = \Theta\left(\sqrt{\frac{1}{N} \log \frac{|\mathcal{F}|}{\delta}}\right)$ .

**Policy Evaluation Error.** The performance of the proposed SPECTRALDICE algorithm is evaluated by the *policy evaluation error*  $\mathcal{E} := \hat{\rho}(\pi) - \rho_{\mathbb{P}}(\pi)$ , which can be further bounded by the following theorem.

**Theorem 4** (Main Theorem). *Suppose Claim 3 holds for the REPLEARN subroutine. Then under Assumptions 1 to 3, with probability at least  $1 - \delta$ , we have*

$$\mathcal{E} \lesssim \frac{1}{1 - \gamma} \sqrt{\frac{\log(1/\delta)}{N}} + \frac{1}{(1 - \gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta/2).$$

*Proof sketch.* We first split  $\mathcal{E}$  into the following terms:

$$\mathcal{E} = \underbrace{\hat{\rho}(\pi) - \bar{\rho}(\pi)}_{\text{statistical}} + \underbrace{\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi)}_{\text{dataset}} + \underbrace{\rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi)}_{\text{representation}},$$

where we introduce an auxiliary problem:

$$\begin{aligned} \bar{\rho}(\pi) = \min_{\theta_Q} \max_{\omega_d} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot | s)} \left[ \hat{\phi}(s, a)^{\top} \theta_Q \right] \right. \\ & + \mathbb{E}_{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a | s), (s', a') \sim \mathbb{P}^{\pi}(\cdot, \cdot | s, a)} \left[ \left( \hat{\mu}^{\pi}(s, a)^{\top} \omega_d \right) \cdot \right. \\ & \left. \left. (r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q) \right] \right\}. \end{aligned}$$

Note that (9) is the empirical estimation of  $\bar{\rho}(\pi)$ , and that  $\bar{\rho}(\pi)$  is (subtly) inequivalent to  $\rho_{\hat{\mathbb{P}}}(\pi)$ —the expectation is still taken over  $(s', a') \sim \mathbb{P}^{\pi}(\cdot, \cdot | s, a)$  rather than  $\hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a) = \langle \hat{\phi}(s, a), \hat{\mu}^{\pi}(\cdot, \cdot) \rangle$ .

Intuitively, the latter two terms are directly related to the representation learning error established in Claim 3, which can actually be bounded as follows:

$$\rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) \lesssim \frac{\gamma}{(1 - \gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta/2),$$

$$\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi) \lesssim \frac{1}{1 - \gamma} \cdot \xi(|\mathcal{F}|, N, \delta/2).$$

On the other hand, the first term is only caused by replacing the expectations with their empirical estimators,

which can be bounded by concentration inequalities as:

$$\hat{\rho}(\pi) - \bar{\rho}(\pi) \lesssim \frac{1}{1 - \gamma} \sqrt{\frac{\log(1/\delta)}{N}}.$$

Plugging these terms back completes the proof.  $\square$

Finally, we conclude that the sample complexity of SPECTRALDICE equipped with either OLS or NCE REPLEARN subroutine is  $\tilde{O}(N^{-1/2})$  (under mild regularity assumptions). Details are deferred to Appendix D.

## 5 EXPERIMENTS

In this section, we present experimental results in both continuous and discrete environments to demonstrate the strength of the proposed SPECTRALDICE algorithm. We also study the impact of hyperparameters, data coverage and the choice of behavior policy on the OPE performance, and illustrate the efficacy of the proposed representation learning method.

The empirical results show that our method outperforms BESTDICE, the state-of-the-art DICE implementation without representation learning, in terms of both the convergence rate and the final prediction error. In comparison to other baselines, SPECTRALDICE achieves comparable performance with higher efficiency in simple environments, and performs significantly better than others in the most challenging environment.

### 5.1 Continuous Environments

**Setting.** We start by comparing SPECTRALDICE with various baseline OPE methods in literature, including BESTDICE (Yang et al., 2020), Fitted Q Evaluation (FQE) (Kostrikov and Nachum, 2020), Model-Based (MB) method (Zhang et al., 2021), Importance Sampling (IS) method (Hanna et al., 2019) and Doubly-Robust (DR) method (Dudík et al., 2011). We follow the experiment protocol in Yang et al. (2020) to evaluate and compare the OPE performances of these algorithms in three continuous MuJoCo environments, namely **Cartpole**, **Reacher** and **Half-Cheetah**, in an increasing order of difficulty. In our implementation, for representation learning, we parameterize each of  $\hat{\phi}$  and  $\hat{\mu}^{\pi}$  with a 2-layer feed-forward neural network. For the OPE step, regularizer is appended to (9), and the estimated policy value is retrieved by  $\hat{\rho}(\pi) = \mathbb{E}_{(s,a) \sim d^{\pi}} [\hat{\mu}^{\pi}(s, a)^{\top} \omega_d \cdot r(s, a)]$  (see Remark 2). Both steps are regarded as stochastic optimization problems, and are solved by stochastic gradient descent and stochastic gradient descent-ascent, respectively. Optimization hyperparameters are selected via grid search. Performance is quantified by *OPE error*  $|\hat{\rho}(\pi) - \rho(\pi)|$ .

**Results.** The OPE performances of different methods in three environments are shown in Figure 1. It is observed that SPECTRALDICE achieves comparable performance in fewer optimization steps as compared to all the other baselines, and further, outperforms

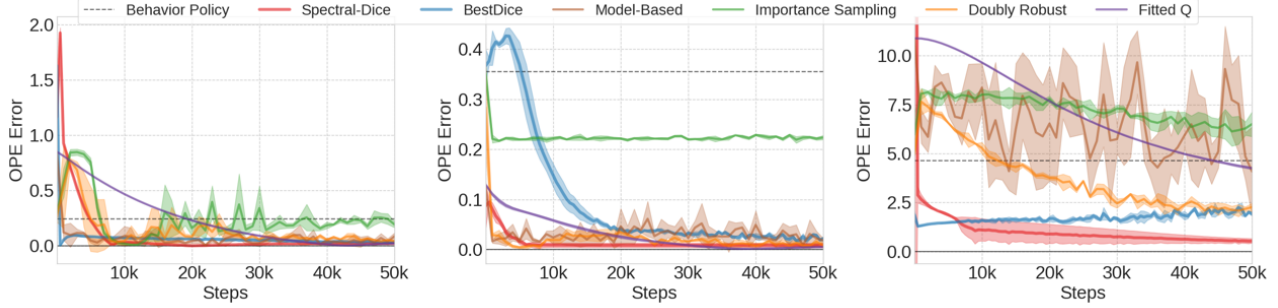


Figure 1: OPE error over the number of training steps in **Cartpole**, **Reacher** and **Half-Cheetah** environments (from left to right). Due to the use of convex-concave formulation, we can see that SPECTRALDICE converges faster and more stably to the target policy with a smaller OPE error in all three environments.

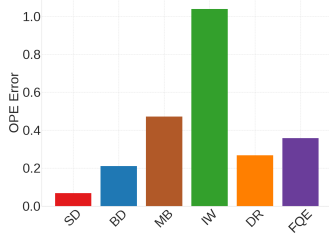


Figure 2: Averaged relative OPE errors over three environments.

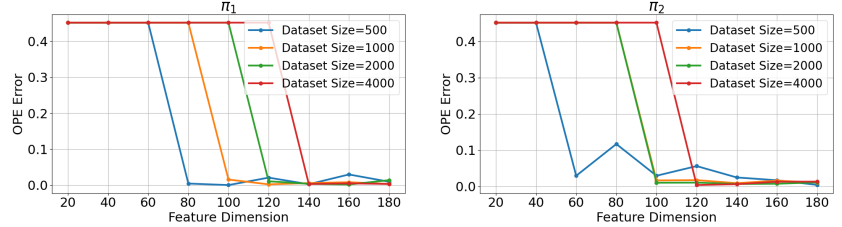


Figure 3: OPE error of SPECTRALDICE in **Four Rooms** with varying behavior policies (“far-away” policy  $\pi_1$  vs. “similar” policy  $\pi_2$ ), dataset sizes and feature dimensions.

them in terms of both convergence rate and final estimation error in the most challenging **Half-Cheetah** environment. Further, although FQE achieves an error close to SPECTRALDICE in simpler environments, its performance significantly degrades when the transition dynamics becomes more complex, demonstrating the importance and power of spectral representation.

Here we also highlight the comparison between two DICE-based methods—SPECTRALDICE (ours) and BESTDICE. All settings showcase the advantage of our primal-dual spectral representation over the generic neural network representation, which justify the argument that, compared to the non-convex non-concave optimization in vanilla DICE, our convex-concave optimization leads to faster convergence and enhanced stability within a wider range of environments.

For a clearer comparison, we further present the averaged relative OPE error across these three environments in Figure 2. Here the *relative OPE error* is defined by  $\frac{|\hat{\rho}(\pi) - \rho(\pi)|}{|\hat{\rho}(\pi_b) - \rho(\pi)|}$ , i.e., OPE error normalized by the value difference between the target and behavior policies. Under this metric, it becomes more evident that our method outperforms all the baselines in terms of estimation accuracy by a large margin.

## 5.2 Discrete Environment

**Setting.** We proceed to test our method in **Four Rooms** (Sutton et al., 1999), a classical discrete environment featuring convenient visualization, to study

the algorithm’s sensitivity for hyperparameters and illustrate the efficacy of representation learning. For representation learning in this tabular MDP, we perform singular value decomposition (SVD) of the matrix  $[\frac{\mathbb{P}^\pi(s', a' | s, a)}{d^\pi(s', a')}]$  (indexed by  $(s, a)$  and  $(s', a')$ ) and select the top  $d$  singular vectors as  $\hat{\phi}(s, a)$  and  $\hat{\mu}^\pi(s', a')$ .

**Sensitivity Study.** We study the algorithm’s sensitivity with respect to behavior policy  $\pi_b$ , dataset size  $N$  and spectral feature dimension  $d$  by examining their impact on the OPE performance. For  $\pi_b$ , we vary between two behavior policies  $\pi_1$  and  $\pi_2$ , where  $\pi_1$  has a larger  $\ell_1$ -distance from the target policy than  $\pi_2$ . The results are shown in Figure 3. It can be observed that the proposed algorithm is always able to achieve low OPE errors with sufficiently large feature dimensions, showcasing its wide applicability under different behavior policies, data availability and hyperparameters.

**Efficacy of Representation Learning.** To give a hint of the efficacy of our representation learning scheme REPLEARN, we visualize in Figure 4 the learned transition kernel  $\hat{\mathbb{P}}$  for a fixed state and all the four actions, where  $\hat{\mathbb{P}}$  is restored from the spectral representation by (5). As shown in the heat map (where darker color indicates higher probability), the REPLEARN algorithm successfully learns a set of primal-dual features that accurately encode the correct transition dynamics.

More experimental details are deferred to Appendix A.



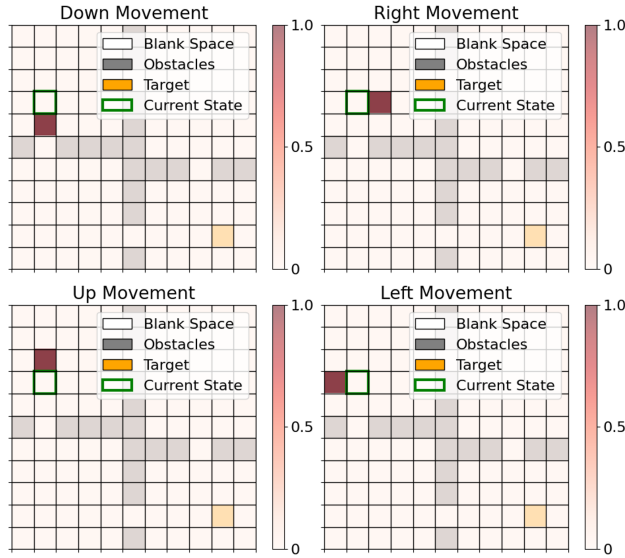


Figure 4: Visualization of the learned transition kernel for a fixed state and all the four actions.

## 6 CONCLUSION

In this paper, to relieve the intrinsic tension between breaking the curse of horizon and overcoming the curse of dimensionality via DICE estimators, we propose a novel primal-dual spectral representation method that establishes linear spectral representations for both the primal variable (*i.e.*,  $Q$ -function) and the dual variable (*i.e.*, stationary distribution correction ratio), which leads to SPECTRALDICE, an efficient and practical OPE algorithm that eliminates the non-convex non-concave saddle-point optimization in DICE and makes efficient use of historical data. The performance of SPECTRALDICE is justified by a theoretical sample complexity guarantee and the empirical outperformance. Future directions include taking one step further to design offline policy optimization methods using primal-dual spectral representations, and applying the algorithm for efficient imitation learning.

## Acknowledgment

This paper is supported in part by NSF ECCS-2328241, NSF CBET-2112085, NSF ECCS-2401390, NSF ECCS-2401391, NSF IIS-2403240, Dolby support, and Schmidt Sciences AI2050 Fellowship.

## References

- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. FLAMBE: Structural complexity and representation learning of low rank MDPs. *Advances in neural information processing systems*, 33:20095–20107, 2020.
- András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71:89–129, 2008.
- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine learning proceedings 1995*, pages 30–37. Elsevier, 1995.
- Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- Justin Boyan and Andrew Moore. Generalization in reinforcement learning: Safely approximating the value function. *Advances in neural information processing systems*, 7, 1994.
- Michel Broniatowski and Amor Keziou. Minimization of  $\phi$ -divergences on sets of signed measures. *Studia Scientiarum Mathematicarum Hungarica*, 43(4):403–442, 2006.
- Jonathan Chang, Kaiwen Wang, Nathan Kallus, and Wen Sun. Learning Bellman complete representations for offline policy evaluation. In *International Conference on Machine Learning*, pages 2938–2971. PMLR, 2022.
- Fengdi Che, Chenjun Xiao, Jincheng Mei, Bo Dai, Ramki Gummadi, Oscar A Ramirez, Christopher K Harris, A Rupam Mahmood, and Dale Schuurmans. Target networks and over-parameterization stabilize off-policy bootstrapping with function approximation. *arXiv preprint arXiv:2405.21043*, 2024.
- Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. Top- $k$  off-policy correction for a reinforce recommender system. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pages 456–464, 2019.
- Bo Dai, Ofir Nachum, Yinlam Chow, Lihong Li, Csaba Szepesvári, and Dale Schuurmans. CoinDICE: Off-policy confidence interval estimation. *Advances in neural information processing systems*, 33:9398–9411, 2020.
- Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. *arXiv preprint arXiv:1103.4601*, 2011.
- Xing Fang, Qichao Zhang, Yinfeng Gao, and Dongbin Zhao. Offline reinforcement learning for autonomous driving with real world driving data. In *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pages 3417–3422. IEEE, 2022.

- Dylan J Foster, Akshay Krishnamurthy, David Simchi-Levi, and Yunzong Xu. Offline reinforcement learning: Fundamental barriers for value function approximation. *arXiv preprint arXiv:2111.10919*, 2021.
- Justin Fu, Mohammad Norouzi, Ofir Nachum, George Tucker, Ziyu Wang, Alexander Novikov, Mengjiao Yang, Michael R Zhang, Yutian Chen, Aviral Kumar, et al. Benchmarks for deep off-policy evaluation. *arXiv preprint arXiv:2103.16596*, 2021.
- Geoffrey J Gordon. Stable function approximation in dynamic programming. In *Machine learning proceedings 1995*, pages 261–268. Elsevier, 1995.
- Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Srinivasan, Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, et al. Evaluating reinforcement learning algorithms in observational health settings. *arXiv preprint arXiv:1805.12298*, 2018.
- Josiah Hanna, Scott Niekum, and Peter Stone. Importance sampling policy evaluation with an estimated behavior policy. In *International Conference on Machine Learning*, pages 2605–2613. PMLR, 2019.
- Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*, 34:5000–5011, 2021.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected works of Wassily Hoeffding*, pages 409–426, 1994.
- Audrey Huang, Jinglin Chen, and Nan Jiang. Reinforcement learning in low-rank MDPs with density features. In *International Conference on Machine Learning*, pages 13710–13752. PMLR, 2023.
- Abhyuday Jagannatha, Philip Thomas, and Hong Yu. Towards high confidence off-policy reinforcement learning for clinical applications. In *CausalML Workshop, ICML*, 2018.
- Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- Haoming Jiang, Bo Dai, Mengjiao Yang, Tuo Zhao, and Wei Wei. Towards automatic evaluation of dialog systems: A model-free off-policy evaluation approach. *arXiv preprint arXiv:2102.10242*, 2021.
- Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International conference on machine learning*, pages 652–661. PMLR, 2016.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are PAC-learnable. In *International Conference on Machine Learning*, pages 1704–1713. PMLR, 2017.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143. PMLR, 2020.
- Gregory Kahn, Adam Villafior, Pieter Abbeel, and Sergey Levine. Composable action-conditioned predictors: Flexible off-policy learning for robot navigation. In *Conference on robot learning*, pages 806–816. PMLR, 2018.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*, pages 651–673. PMLR, 2018.
- Nathan Kallus and Masatoshi Uehara. Double reinforcement learning for efficient off-policy evaluation in markov decision processes. *Journal of Machine Learning Research*, 21(167):1–63, 2020.
- Pulkit Katdare, Nan Jiang, and Katherine Rose Driggs-Campbell. Marginalized importance sampling for off-environment policy evaluation. In *Conference on Robot Learning*, pages 3778–3788. PMLR, 2023.
- Ilya Kostrikov and Ofir Nachum. Statistical bootstrapping for uncertainty estimation in off-policy evaluation. *arXiv preprint arXiv:2007.13609*, 2020.
- Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In *International Conference on Machine Learning*, pages 3703–3712. PMLR, 2019.
- Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-bandit-based news article recommendation algorithms. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 297–306, 2011.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. *Advances in neural information processing systems*, 31, 2018.
- Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *AAMAS*, volume 1077, 2014.
- Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9(5), 2008.

- Ofir Nachum and Bo Dai. Reinforcement learning via Fenchel-Rockafellar duality. *arXiv preprint arXiv:2001.01866*, 2020.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in neural information processing systems*, 32, 2019a.
- Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. AlgaeDICE: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019b.
- Chengzhuo Ni, Anru R Zhang, Yaqi Duan, and Mengdi Wang. Learning good state and action representations via tensor decomposition. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pages 1682–1687. IEEE, 2021.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Shuang Qiu, Lingxiao Wang, Chenjia Bai, Zhuoran Yang, and Zhaoran Wang. Contrastive UCB: Provably efficient contrastive self-supervised learning in online reinforcement learning. In *International Conference on Machine Learning*, pages 18168–18210. PMLR, 2022.
- Tongzheng Ren, Chenjun Xiao, Tianjun Zhang, Na Li, Zhaoran Wang, Sujay Sanghavi, Dale Schuurmans, and Bo Dai. Latent variable representation for reinforcement learning. *arXiv preprint arXiv:2212.08765*, 2022a.
- Tongzheng Ren, Tianjun Zhang, Lisa Lee, Joseph E Gonzalez, Dale Schuurmans, and Bo Dai. Spectral decomposition representation for reinforcement learning. *arXiv preprint arXiv:2208.09515*, 2022b.
- Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical exploration for representation learning. In *Uncertainty in Artificial Intelligence*, pages 1686–1696. PMLR, 2022c.
- Tongzheng Ren, Zhaolin Ren, Na Li, and Bo Dai. Stochastic nonlinear control via finite-dimensional spectral dynamic embedding. *arXiv preprint arXiv:2304.03907*, 2023.
- Tianyu Shi, Dong Chen, Kaian Chen, and Zhaojian Li. Offline reinforcement learning for autonomous driving with safety and exploration enhancement. *arXiv preprint arXiv:2110.07067*, 2021.
- Dmitry Shribak, Chen-Xiao Gao, Yitong Li, Chenjun Xiao, and Bo Dai. Diffusion spectral representation for reinforcement learning. *arXiv preprint arXiv:2406.16121*, 2024.
- Ahmad Slim, Husain Al Yusuf, Nadine Abbas, Chaouki T Abdallah, Gregory L Heileman, and Ameer Slim. A Markov decision processes modeling for curricular analytics. In *2021 20th IEEE international conference on machine learning and applications (ICMLA)*, pages 415–421. IEEE, 2021.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 2139–2148. PMLR, 2016.
- John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. *Advances in neural information processing systems*, 9, 1996.
- Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under partial coverage. *arXiv preprint arXiv:2107.06226*, 2021.
- Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and  $Q$ -function learning for off-policy evaluation. In *International Conference on Machine Learning*, pages 9659–9668. PMLR, 2020.
- Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline RL in low-rank MDPs. *arXiv preprint arXiv:2110.04652*, 2021.
- Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. *Advances in neural information processing systems*, 32, 2019.
- Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR, 2020.
- Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy evaluation via the regularized Lagrangian. *Advances in Neural Information Processing Systems*, 33:6551–6561, 2020.
- Mengjiao Yang, Bo Dai, Ofir Nachum, George Tucker, and Dale Schuurmans. Offline policy selection under uncertainty. In *International Conference on Artificial Intelligence and Statistics*, pages 4376–4396. PMLR, 2022.
- Hengshuai Yao, Csaba Szepesvári, Bernardo Avila Pires, and Xinhua Zhang. Pseudo-MDPs and factored linear action models. In *2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL)*, pages 1–9. IEEE, 2014.

- Wenhao Zhan, Baihe Huang, Audrey Huang, Nan Jiang, and Jason Lee. Offline reinforcement learning with realizability and single-policy concentrability. In *Conference on Learning Theory*, pages 2730–2775. PMLR, 2022.
- Hongming Zhang, Tongzheng Ren, Chenjun Xiao, Dale Schuurmans, and Bo Dai. Provable representation with efficient planning for partially observable reinforcement learning. *arXiv preprint arXiv:2311.12244*, 2023.
- Michael R Zhang, Tom Le Paine, Ofir Nachum, Cosmin Paduraru, George Tucker, Ziyu Wang, and Mohammad Norouzi. Autoregressive dynamics models for offline policy evaluation and optimization. *arXiv preprint arXiv:2104.13877*, 2021.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020.
- Tianjun Zhang, Tongzheng Ren, Mengjiao Yang, Joseph Gonzalez, Dale Schuurmans, and Bo Dai. Making linear MDPs practical via contrastive representation learning. In *International Conference on Machine Learning*, pages 26447–26466. PMLR, 2022a.
- Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Alekh Agarwal, and Wen Sun. Efficient reinforcement learning in block MDPs: A model-free representation learning approach. In *International Conference on Machine Learning*, pages 26517–26547. PMLR, 2022b.



## Checklist

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

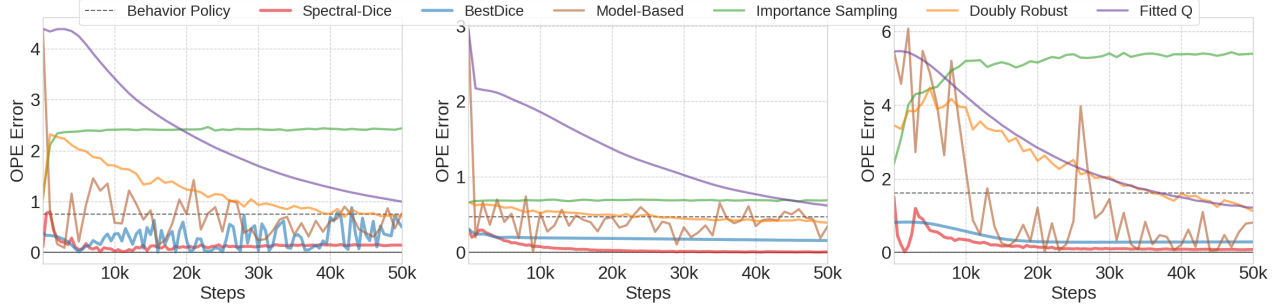


Figure 5: OPE error over the number of training steps in **Walker2d**, **Hopper** and **Ant** (from left to right).

## A More Experimental Results

**Additional Experiments.** We evaluate the OPE performance of the proposed SPECTRALDICE algorithm and the aforementioned baselines (see Section 5.1) in three additional environments, namely **Walker2d**, **Hopper** and **Ant**, the results of which are shown in Figure 5. These additional experiments further justify that our algorithm outperforms all the other baselines in a consistent and robust way, enjoying both a faster convergence rate and a smaller OPE error. These additional experimental results further confirm the superiority of SPECTRALDICE.

**Average Error.** The average loss of all the methods across three environments (**Cartpole**, **Reacher** and **Half-Cheetah**) can be found below in Table 1. Experiments are repeated using three random seeds.

Methods	SPECTRALDICE	BESTDICE	MB	IS	DR	FQE
Average Loss	$0.0781 \pm 0.0070$	$0.2147 \pm 0.0077$	$0.5280 \pm 0.0442$	$0.9734 \pm 0.0538$	$0.2559 \pm 0.0088$	$0.3474 \pm 0.0084$

Table 1: Overall results across different methods.

**Learning Efficiency.** We further evaluate the efficiency of the proposed SPECTRALDICE algorithm by comparing its running time against the BESTDICE baseline in multiple environments, the results of which are shown in Table Appendix A. Here the training process is stopped after the evaluation loss drops below a preset threshold for 5 straight test epochs. It can be observed that Stage 2 (line 2 in Algorithm 1) indeed features accelerated updates, and even with the additional overhead induced by Stage 1 (line 1 in Algorithm 1), our SPECTRALDICE algorithm becomes advantageous over the baseline for more challenging tasks.

Intuitively, despite the additional overhead induced by representation learning in Stage 1, the optimization in Stage 2 is more efficient than its vanilla counterpart, since the updates are only performed with respect to cofactors  $\theta_Q$  and  $\omega_d$  that lie in a low-dimensional space. Consequently, Stage 2 by itself is expected to outperform vanilla DICE in terms of both memory use and running time.

Method	Cartpole	Reacher	Half-Cheetah
SPECTRALDICE, Stage 1 (s)	715.9	876.0	913.0
SPECTRALDICE, Stage 2 (s)	167.5	140.8	357.1
BESTDICE (s)	482.0	1676.5	1989.0

Table 2: Training time comparison between SpectralDice and BestDice

**Implementation Details.** For the baseline algorithms, we follow the implementation of BESTDICE in Yang et al. (2020) and the implementations of FQE, MB, IS, DR in Fu et al. (2021). The optimization hyperparameters including learning rate, optimizer parameter, network architecture, batch size, *etc.*, are selected via grid search. All the experiments were conducted using V100 GPUs on a multi-node cluster.

For the continuous environments, the target policy is obtained using deep reinforcement learning agents (Deep Q-Network (DQN) agent for **Cartpole**, and Soft Actor-Critic (SAC) agents for all the other environments). The behavior policy is then obtained by sampling from a Gaussian distribution centered at the mean action of the target policy, where the variance of the Gaussian distribution can be adjusted to get behavior policies at different distances from the target policy. To build the offline dataset, we collect 400 trajectories using the behavior policy, where each trajectory is truncated to 250 steps.

The source code is available at [https://anonymous.4open.science/r/spectral\\_dice-720A](https://anonymous.4open.science/r/spectral_dice-720A).

## B Primal-Dual Spectral Representation

In this appendix, we present the key properties of the proposed primal-dual spectral representation with proofs, as well as a brief discussion on why the spectral representation of that specific form is preferable.

### B.1 DICE Estimator with Primal-Dual Spectral Representation

We first present the proof of Corollary 2 that is already stated in the main text.

**Corollary 2.** *With primal-dual spectral representation (6) where  $q(\cdot) \equiv d^{\pi_b}(\cdot)$ , under Assumption 1, we have*

$$\begin{aligned} \rho(\pi) = \min_{\theta_Q} \max_{\omega_d} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [\phi(s, a)^\top \theta_Q] \right. \\ & \left. + \mathbb{E}_{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [(\mu^\pi(s, a)^\top \omega_d) (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q)] \right\}. \end{aligned}$$

*Proof of Corollary 2.* Recall the primal-dual LP formulation of policy evaluation stated in (3), which can be equivalently rewritten using the primal-dual spectral representation (6) as follows:

$$\rho(\pi) = \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int d^\pi(s, a) \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] ds da \right\} \quad (11a)$$

$$= \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int q(s) \pi_b(a|s) \cdot \frac{d^\pi(s, a)}{q(s) \pi_b(a|s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] ds da \right\} \quad (11b)$$

$$= \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \mathbb{E}_{s \sim q(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[ \frac{d^\pi(s, a)}{q(s) \pi_b(a|s)} (r(s, a) + \gamma Q(s', a') - Q(s, a)) \right] \right\} \quad (11c)$$

$$= \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [\phi(s, a)^\top \theta_Q] + \mathbb{E}_{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [\mu^\pi(s, a)^\top \omega_d (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q)] \right\}, \quad (11d)$$

where in (11b) we perform the IS-style change-of-variable used in DICE estimators (see (4)); in (11d) we plug in the primal-dual spectral representation of  $Q^\pi$  and  $d^\pi$  stated in (6), as well as the fact that  $q(\cdot) \equiv d^{\pi_b}(\cdot)$ .  $\square$

### B.2 Failure of the Naive Spectral Representation

In Section 3.1, it is mentioned that directly applying the naive spectral representation (5) proposed in Ren et al. (2022b) induces a complicated representation for  $\zeta(\cdot, \cdot)$ , which in turn leads to an intractable optimization (3) for the computation of the DICE estimator. The above point is further elaborated here in a formal way.

Note that, in Lemma 1, the linear representation of  $Q^\pi$  only builds upon the low-rank MDP assumption, and therefore it still holds with the naive spectral representation (5). On the other hand, it can be checked that

$$d^\pi(s, a) = \underbrace{\left\langle q(s) \pi(a|s) \mu(s), (1 - \gamma) \omega_0 + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \phi(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a} \right\rangle}_{\omega_d^\pi}, \quad (12)$$

which can be obtained by plugging the relation  $\pi(a|s) \mu(s) = \pi_b(a|s) \mu^\pi(s, a)$  into the linear representation of  $d^\pi(\cdot, \cdot)$  to eliminate  $\mu^\pi$  from the representation. Consequently, the LP formulation (11) becomes

$$\begin{aligned} \rho(\pi) = \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int d^\pi(s, a) \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] ds da \right\} \\ = \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int q(s) \pi_b(a|s) \cdot \frac{\pi(a|s)}{\pi_b(a|s)} \frac{d^\pi(s, a)}{q(s) \pi(a|s)} \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] ds da \right\} \\ = \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \mathbb{E}_{s \sim q(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} \frac{d^\pi(s, a)}{q(s) \pi(a|s)} (r(s, a) + \gamma Q(s', a') - Q(s, a)) \right] \right\} \\ = \min_{\theta_Q} \max_{\omega_d} & \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [\phi(s, a)^\top \theta_Q] + \mathbb{E}_{s \sim q(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} (\mu(s)^\top \omega_d) (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q) \right] \right\}, \end{aligned}$$

which involves an unknown policy ratio  $\frac{\pi(a|s)}{\pi_b(a|s)}$  when the behavior policy  $\pi_b$  is unknown, and is thus intractable.

The above failed attempt implies that the policy ratio should be “absorbed” into the representation to be implicitly learned during representation learning, which exactly inspires the primal-dual spectral representation (6).

### B.3 Solving the Minimax Problem via Regularization

It is known that directly solving (9) leads to potential numerical instability issues due to the objective's linearity in  $\theta_Q$  and  $\omega_d$  (Nachum et al., 2019b). Fortunately, it is shown in Yang et al. (2020) that certain regularization leads to strictly concave inner maximization while keeping the optimal solution  $\omega_d^*$  unbiased. Specifically, in practical implementation we may append the following regularizer to the objective in (10):

$$\rho_{\text{reg}}(\pi) = \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} [\phi(s, a)^\top \theta_Q] + \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), \ a \sim \pi_b(a|s), \\ s' \sim \mathbb{P}(\cdot|s, a), \ a' \sim \pi(\cdot|s')}} \left[ (\mu^\pi(s, a)^\top \omega_d) \cdot \right. \right. \\ \left. \left. (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q) \right] - \lambda \mathbb{E}_{(s, a) \sim \mathcal{D}} [f(\hat{\mu}^\pi(s, a)^\top \omega_d)] \right\}. \quad (13)$$

Here  $f$  is a differentiable convex function with closed and convex Fenchel conjugate  $f_*$  (see Appendix E.1), and  $\lambda > 0$  is a tunable constant that controls the magnitude of regularization. It is evident that the regularized objective is concave in  $\omega_d$ , which facilitates the inner maximization. What's more, it has also been shown that such regularization does not alter the optimal solution  $\omega_d^*$ , as summarized in the following lemma.

**Lemma 5** (Nachum et al. (2019b); Yang et al. (2020)). *The solution  $(\theta_Q^{\text{reg},*}, \omega_d^{\text{reg},*})$  to (13) is characterized by:*

$$\begin{aligned} \phi(s, a)^\top \theta_Q^{\text{reg},*} &= \phi(s, a)^\top \theta_Q^* - \lambda (\mathcal{I} - \mathcal{P}^\pi)^{-1} f' \left( \frac{d^\pi(s, a)}{d^{\pi_b}(s, a)} \right), \\ \mu^\pi(s, a)^\top \omega_d^* &= \mu^\pi(s, a)^\top \omega_d^{\text{reg},*}, \\ \rho_{\text{reg}}(\pi) &= \rho(\pi) - \lambda D_f(d^\pi \| d^{\pi_b}), \end{aligned}$$

where  $(\theta_Q^*, \omega_d^*)$  is the solution to (10).

We emphasize that the regularized problem is unbiased only in the sense that  $\omega_d^{\text{reg},*} = \omega_d^*$ . Therefore, in general we need to plug  $\omega_d^{\text{reg},*}$  back into (10) and solve the outer minimization again to recover  $\theta_Q^*$ . Nevertheless, when  $\lambda$  is sufficiently small, we shall regard  $\theta_Q^{\text{reg},*} \approx \theta_Q^*$  to relieve the additional computational burden.

In practice, we can only solve the empirical version of (13), i.e.,

$$\rho_{\text{reg}}(\pi) = \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \widehat{\mathbb{E}}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} [\phi(s, a)^\top \theta_Q] + \widehat{\mathbb{E}}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} \left[ (\mu^\pi(s, a)^\top \omega_d) \cdot \right. \right. \\ \left. \left. (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q) \right] - \lambda \widehat{\mathbb{E}}_{(s, a) \sim \mathcal{D}} [f(\hat{\mu}^\pi(s, a)^\top \omega_d)] \right\}.$$

## C Representation Learning Methods and Their Error Bounds

In this appendix, we introduce two candidate methods—*ordinary least squares (OLS)* and *noise-contrastive estimation (NCE)*—that can be used as the REPLEARN subroutine. Further, we also provide their representation learning error bounds in the form of Claim 3, which is restated here for readers' convenience:

**Claim 3.** *With probability at least  $1 - \delta$ , the representation learning error of REPLEARN( $\mathcal{F}, \mathcal{D}, \pi$ ) is bounded by*

$$\mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[ \left\| \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) - \mathbb{P}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \leq \xi(|\mathcal{F}|, N, \delta),$$

where  $\hat{\mathbb{P}}^\pi(s', a' | s, a) := d^{\pi_b}(s', a') \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s')$ , and  $N$  is the number of samples in  $\mathcal{D}$ .

It should be emphasized that the two methods discussed here are not the only candidates for REPLEARN. Rather, any representation learning method that comes with a learning error bound in the required form is applicable, without any further requirements on the learning mechanism.

### C.1 Ordinary Least Squares (OLS)

**Method.** Inspired by Ren et al. (2022b), the objective of OLS can be constructed as follows. Denote by  $\mathbb{Q}^\pi(s', a', s, a) := d^{\pi_b}(s, a) \mathbb{P}^\pi(s', a' | s, a)$  the joint distribution of state-action transitions under behavior policy  $\pi_b$ . Then we plug  $\mathbb{Q}^\pi$  into (6) and rearrange the terms to obtain

$$\frac{\mathbb{Q}^\pi(s', a', s, a)}{\sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')}} = \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \phi(s, a)^\top \mu^\pi(s', a').$$

Therefore, we propose to optimize over the following OLS objective:

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \int \left( \frac{\mathbb{Q}^\pi(s', a', s, a)}{\sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')}} - \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right)^2 ds da ds' da'$$



$$\begin{aligned}
 &= \min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \left\{ \int \frac{\mathbb{Q}^\pi(s', a', s, a)^2}{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} ds da ds' da' - 2 \mathbb{E}_{(s, a) \sim d^{\pi_b}, (s', a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)} \left[ \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right] \right. \\
 &\quad \left. + \mathbb{E}_{(s, a) \sim d^{\pi_b}, (s', a') \sim d^{\pi_b}} \left[ (\hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a'))^2 \right] \right\},
 \end{aligned}$$

Note that the first term  $\int \frac{\mathbb{Q}^\pi(s', a', s, a)^2}{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} ds da ds' da'$  is a constant that can be omitted in optimization, while the second and third terms can be effectively approximated by sampling from the dataset  $\mathcal{D}$  and the target policy  $\pi$ . Therefore, in practice we learn  $(\hat{\phi}, \hat{\mu}^\pi)$  by solving the following optimization:

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \left\{ \widehat{\mathbb{E}}_{(s, a) \sim d^{\pi_b}, (\tilde{s}', \tilde{a}') \sim d^{\pi_b}} \left[ (\hat{\phi}(s, a)^\top \hat{\mu}^\pi(\tilde{s}', \tilde{a}'))^2 \right] - 2 \widehat{\mathbb{E}}_{(s, a) \sim d^{\pi_b}, (s', a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)} \left[ \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right] \right\}, \quad (14)$$

where the expectations are replaced by their empirical estimations using data sampled from  $\mathcal{D}$ .

**Error Bound.** We proceed to show the representation learning error bound for the OLS method, which requires the following regularity assumption on the transition kernel  $\mathbb{P}^\pi$  and the occupancy measure  $d^{\pi_b}$ .

**Assumption 4** (regularity for OLS). (1) lower-bounded transition kernel:  $\mathbb{P}^\pi(s', a' | s, a) \geq \frac{1}{C_P} > 0$ ,  $\forall s, a, s', a'$ ; (2) effective behavior policy coverage:  $\frac{d^{\pi_b}(s, a)}{d^{\pi_b}(s', a')} \leq C_{\text{cov}}$ ,  $\forall s, a, s', a'$ .

We point out that the major rationale behind these mild assumptions is to rule out the cases where certain transitions are scarcely sampled due to the singularity in transition kernel or behavior policy.

**Theorem 6** (OLS learning error). *Under Assumptions 1 to 3 and the additional Assumption 4 for regularity, let  $(\hat{\phi}, \hat{\mu}^\pi)$  be the solution to (14), and set  $\hat{\mathbb{P}}^\pi(s', a' | s, a) := d^{\pi_b}(s', a') \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s')$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have*

$$\mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[ \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \leq \sqrt{C_P C_{\text{reg}}} \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}},$$

where  $C_{\text{reg}} = \frac{4}{3} \sqrt{C_{\text{cov}}} + 8C_{\text{cov}}$  is a universal constant determined by the PAC bound for OLS.

*Proof.* We would like to apply the fast-rate PAC bound for OLS regression (Lemma 18). For the sake of clarity, we explicitly define the family of candidate regression functions as

$$\tilde{\mathcal{F}} := \left\{ f : (s, a, s', a') \mapsto \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \phi(s, a)^\top \mu^\pi(s', a') \mid (\phi, \mu^\pi) \in \mathcal{F} \right\}.$$

It is evident that any  $f \in \tilde{\mathcal{F}}$  is bounded as follows:

$$0 \leq f(s, a, s', a') = \sqrt{\frac{d^{\pi_b}(s, a)}{d^{\pi_b}(s', a')}} \tilde{\mathbb{P}}^\pi(s', a' | s, a) \leq \sqrt{C_{\text{cov}}},$$

where we use the fact that  $\langle \hat{\phi}(s, a), d^{\pi_b}(s', a') \hat{\mu}^\pi(s', a') \rangle$  is always some valid transition kernel  $\tilde{\mathbb{P}}^\pi$  (by Assumption 3), and the additional regularity assumption (Assumption 4). Further, since the family  $\tilde{\mathcal{F}}$  is realizable (by Assumption 3), there exists an optimal  $f^* \in \tilde{\mathcal{F}}$  such that

$$f^*(s, a, s', a') = \frac{\mathbb{Q}^\pi(s', a', s, a)}{\sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')}}.$$

As  $f(s, a, s', a'), f^*(s, a, s', a') \in [0, \sqrt{C_{\text{cov}}}]$ , we deduce from Lemma 18 that, with probability at least  $1 - \delta$ ,

$$\int \left( f^*(s, a, s', a') - \hat{f}(s, a, s', a') \right)^2 ds da ds' da' \leq C_{\text{reg}} \cdot \frac{\log(|\mathcal{F}|/\delta)}{N}, \quad (15)$$

where  $C_{\text{reg}} := \frac{4}{3} \sqrt{C_{\text{cov}}} + 8C_{\text{cov}}$ , and  $\hat{f}(s, a, s', a') := \sqrt{d^{\pi_b}(s, a) d^{\pi_b}(s', a')} \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a')$ . Consequently,

$$\begin{aligned}
 &\mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[ \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \\
 &= \int d^{\pi_b}(s, a) \left| \mathbb{P}^\pi(s', a' | s, a) - \hat{\mathbb{P}}^\pi(s', a' | s, a) \right| ds da ds' da' \quad (16a)
 \end{aligned}$$

$$= \int \left| \mathbb{Q}^\pi(s', a', s, a) - \hat{\mathbb{Q}}^\pi(s', a', s, a) \right| ds da ds' da' \quad (16b)$$

$$\leq \sqrt{\int \left( \sqrt{\mathbb{Q}^\pi(s', a', s, a)} - \frac{\hat{\mathbb{Q}}^\pi(s', a', s, a)}{\sqrt{\mathbb{Q}^\pi(s', a', s, a)}} \right)^2 ds da ds' da' \cdot \int \mathbb{Q}^\pi(s', a', s', a') ds da ds' da'} \quad (16c)$$

$$= \sqrt{\int \frac{d^{\pi_b}(s', a')}{\mathbb{P}^\pi(s', a' | s, a)} \left( f^*(s, a, s', a') - \hat{f}(s, a, s', a') \right)^2 ds da ds' da'} \quad (16d)$$

$$\leq \sqrt{\max_{s,a,s',a'} \left\{ \frac{d^{\pi_b}(s',a')}{\mathbb{P}^\pi(s',a'|s,a)} \right\}} \cdot \sqrt{C_{\text{reg}} \cdot \frac{\log(|\mathcal{F}|/\delta)}{N}} \quad (16e)$$

$$\leq \sqrt{C_{\mathbb{P}} C_{\text{reg}}} \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}, \quad (16f)$$

where in (16b) we use the definition of  $\mathbb{Q}^\pi$ , and define  $\hat{\mathbb{Q}}^\pi := d_{\mathbb{P}}^{\pi_b}(s,a) \hat{\mathbb{P}}^\pi(s',a'|s,a)$ ; in (16c) we use Cauchy-Schwartz inequality; in (16d) we use the definition of  $\hat{f}$  and  $f^*$ ; in (16e) we plug in the PAC bound (15); in (16f) we use Assumption 4 to bound the coefficient. This completes the proof.  $\square$

## C.2 Noise-Contrastive Learning (NCE)

**Method.** NCE is a widely used method for contrastive representation learning in RL (Zhang et al., 2022a; Qiu et al., 2022). To learn  $(\hat{\phi}, \hat{\mu}^\pi)$ , we consider a binary contrastive learning objective (Qiu et al., 2022):

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \mathbb{E}_{(s,a) \sim d^{\pi_b}} \left[ \mathbb{E}_{(s',a') \sim \mathbb{P}^\pi(\cdot, \cdot | s,a)} \left[ \log \left( 1 + \frac{1}{\hat{\phi}(s,a)^\top \hat{\mu}^\pi(s',a')} \right) \right] + \mathbb{E}_{(s',a') \sim P_{\text{neg}}} \left[ \log \left( 1 + \hat{\phi}(s,a)^\top \hat{\mu}^\pi(s',a') \right) \right] \right], \quad (17)$$

where  $P_{\text{neg}}$  is a negative sampling distribution that will be specified with justification later. We highlight that the above objective implicitly guarantees an equal number of positive and negative samples.

The following derivations follow a similar pathway as those in Qiu et al. (2022). For notational consistency that facilitates the application of known results, we introduce the following auxiliary notations. Define

$$\tilde{\mathcal{F}} := \{f : (s,a,s',a') \mapsto \phi(s,a)^\top \mu^\pi(s',a') \mid (\phi, \mu^\pi) \in \mathcal{F}\}.$$

For clarity, we augment the sampled transitions to include a label  $y$  indicating whether the sample is positive ( $y = 1$ ) or negative ( $y = 0$ ). Formally, given a dataset  $\mathcal{D} = \{(s_i, a_i, s'_i, a'_i) \mid i \in [N]\}$  of positive transitions, we randomly sample  $N$  negative transitions  $(\tilde{s}_i, \tilde{a}_i) \sim P_{\text{neg}}$  ( $i \in [N]$ , i.i.d.), and define the augmented dataset

$$\tilde{\mathcal{D}} := \{(s_i, a_i, s'_i, a'_i, 1), (s_i, a_i, \tilde{s}_i, \tilde{a}_i, 0) \mid i \in [N]\}.$$

In this way, the NCE objective (17) can be equivalently rewritten (in MLE format) as

$$\max_{f \in \tilde{\mathcal{F}}} \mathbb{E}_{(s,a,s',a',y) \sim d^{\tilde{\mathcal{D}}}} [\log \psi_f(s,a,s',a',y)], \quad (18)$$

where the likelihood function  $\psi_f$  is defined by

$$\psi_f(s,a,s',a',y) := \left( \frac{f(s,a,s',a')}{1 + f(s,a,s',a')} \right)^y \cdot \left( \frac{1}{1 + f(s,a,s',a')} \right)^{1-y}.$$

We point out that  $\psi_f(s,a,s',a', \cdot) \in \Delta(\mathcal{Y})$  for any  $(s,a,s',a')$ , where  $\mathcal{Y} := \{0,1\}$ . In fact, given  $f^*$  that optimizes the unconstrained non-empirical version of (18),  $\psi_{f^*}$  can be interpreted as the probability of obtaining label  $y$  given  $(s,a,s',a')$ , as summarized in the following lemma that is similar to Lemma C.1 in Qiu et al. (2022).

**Lemma 7** (non-empirical solution to NCE). *The optimal solution  $f^* := \max_f \mathbb{E}_{(s,a,s',a',y) \sim d^{\tilde{\mathcal{D}}}} [\log \psi_f(s,a,s',a',y)]$  to the unconstrained non-empirical version of (18) is characterized by*

$$f^*(s,a,s',a') = \frac{\mathbb{P}^\pi(s',a'|s,a)}{P_{\text{neg}}(s',a')}.$$

*Proof.* Note that the objective can be rewritten as

$$\begin{aligned} & \mathbb{E}_{(s,a,s',a',y) \sim d^{\tilde{\mathcal{D}}}} [\log \psi_f(s,a,s',a',y)] \\ &= \int d^{\tilde{\mathcal{D}}}(s,a,s',a') \left( \sum_{y \in \mathcal{Y}} \Pr(y|s,a,s',a') \log \psi_f(s,a,s',a',y) \right) ds da ds' da' \\ &= - \int d^{\tilde{\mathcal{D}}}(s,a,s',a') \cdot H(\Pr(y|s,a,s',a'); \psi_f(s,a,s',a',y)) ds da ds' da'. \end{aligned}$$

Here  $H(\cdot; \cdot)$  denotes the cross entropy between distributions, which, by Gibbs' inequality, is minimized only when

$$\Pr(y|s,a,s',a') = \psi_{f^*}(s,a,s',a',y) = \left( \frac{f^*(s,a,s',a')}{1 + f^*(s,a,s',a')} \right)^y \cdot \left( \frac{1}{1 + f^*(s,a,s',a')} \right)^{1-y}. \quad (19)$$

On the other hand, Bayes' rule states that (note that  $\Pr(y|s,a) = \frac{1}{2}$ ,  $\forall y \in \mathcal{Y}$ ):

$$\Pr(y = 1|s,a,s',a') = \frac{\Pr(s',a'|s,a,y=1)\Pr(y=1|s,a)}{\sum_{y \in \mathcal{Y}} \Pr(s',a'|s,a,y)\Pr(y|s,a)} = \frac{\mathbb{P}^\pi(s',a'|s,a)}{P_{\text{neg}}(s',a') + \mathbb{P}^\pi(s',a'|s,a)}. \quad (20)$$

Comparing (19) and (20) gives

$$\frac{f^*(s, a, s', a')}{1 + f^*(s, a, s', a')} = \frac{\mathbb{P}^\pi(s', a' | s, a)}{P_{\text{neg}}(s', a') + \mathbb{P}^\pi(s', a' | s, a)} \implies f^*(s, a, s', a') = \frac{\mathbb{P}^\pi(s', a' | s, a)}{P_{\text{neg}}(s', a')}.$$

This completes the proof.  $\square$

*Remark 4.* For conciseness, here we slightly abuse the notation  $\Pr(\cdot)$  to denote the distribution (density or mass) of joint and conditional distributions involving random variables  $(s, a, s', a', y) \sim d^{\tilde{\mathcal{D}}}$ . Specifically, we write  $\Pr(\cdots, x, \cdots)$  to indicate an arbitrary value  $x$  taken by the random variable, and we also write  $\Pr(\cdots, x = x_0, \cdots)$  to emphasize the specific value  $x_0$  taken by that random variable.

Lemma 7 is important in that it echoes the form of primal-dual spectral representation in (6). Specifically, we shall take  $P_{\text{neg}}(\cdot, \cdot) \equiv d^{\pi_b}(\cdot, \cdot)$  for an exact match, which is also implementable using offline data since  $d^{\pi_b}$  can be effectively approximated by sampling the trajectories. We will stick to this choice of  $P_{\text{neg}}$  from now on.

**Error Bound.** We proceed to show the representation learning error bound for the NCE method, which requires the following regularity assumption on the negative sampling distribution  $P_{\text{neg}}$ , or equivalently, as per the choice above, the state-action occupancy measure  $d^{\pi_b}(\cdot, \cdot)$  for the behavior policy  $\pi_b$ .

**Assumption 5** (regularity for NCE).  $d^{\pi_b}(s, a) \geq \frac{1}{C_d} > 0, \forall s, a$ .

We point out that Assumption 5 is a standard assumption for the negative sampling distribution (Qiu et al., 2022), aiming at eliminating the cases where certain transitions are scarcely drawn as negative samples and thus obstruct efficient representation learning for those cases. The assumption is also slightly stronger than the effective behavior policy coverage assumption required by the OLS method (see Assumption 4).

**Theorem 8** (NCE learning error). *Under Assumptions 1 to 3 and the additional Assumption 5 for regularity, let  $(\hat{\phi}, \hat{\mu}^\pi)$  be the solution to (17) with  $P_{\text{neg}}(\cdot, \cdot) \equiv d^{\pi_b}(\cdot, \cdot)$ , and set  $\hat{\mathbb{P}}^\pi(s', a' | s, a) := d^{\pi_b}(s', a') \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s')$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have*

$$\mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[ \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \leq 2\sqrt{2}(1 + C_d) \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}.$$

The proof of Theorem 8 largely follows the same pathway and techniques established in Qiu et al. (2022). Nevertheless, our proof is less technically involved since the offline non-episodic setting significantly weakens the correlation between samples. For the sake of completeness, we restate the complete proof below.

*Proof.* We start by observing  $\Pr(y, s', a' | s, a) := \Pr(y | s, a, s', a') \Pr(s', a' | s, a)$ , where  $\Pr(s', a' | s, a)$  can in turn be calculated using Bayes' rule as follows:

$$\begin{aligned} \Pr(s', a' | s, a) &= \Pr(s', a' | s, a, y = 0) \Pr(y = 0 | s, a) + \Pr(s', a' | s, a, y = 1) \Pr(y = 1 | s, a) \\ &= \frac{1}{2} (\mathbb{P}^\pi(s', a' | s, a) + P_{\text{neg}}(s', a')). \end{aligned} \quad (21)$$

Here we use the fact that the data distribution  $d^{\tilde{\mathcal{D}}}$  implicitly assigns an equal number of labels as  $y = 0$  and  $y = 1$  by the design of NCE objective (17). Since  $\Pr(s', a' | s, a)$  is a constant that is independent from  $f$ , we can further rewrite the NCE objective to be

$$\arg \max_{f \in \tilde{\mathcal{F}}} \left\{ \mathbb{E}_{(s, a, s', a', y) \in \tilde{\mathcal{D}}} [\log \Pr_f(y | s, a, s', a')] \right\} = \arg \max_{f \in \tilde{\mathcal{F}}} \left\{ \mathbb{E}_{(s, a, s', a', y) \in \tilde{\mathcal{D}}} [\log \Pr_f(y, s', a' | s, a)] \right\}, \quad (22)$$

where we define the shorthand notations

$$\begin{aligned} \Pr_f(y | s, a, s', a') &:= \left( \frac{f(s, a, s', a')}{1 + f(s, a, s', a')} \right)^y \cdot \left( \frac{1}{1 + f(s, a, s', a')} \right)^{1-y}, \\ \Pr_f(y, s', a' | s, a) &:= \left( \frac{f(s, a, s', a') \Pr(s', a' | s, a)}{1 + f(s, a, s', a')} \right)^y \cdot \left( \frac{\Pr(s', a' | s, a)}{1 + f(s, a, s', a')} \right)^{1-y} \end{aligned}$$

for any  $f \in \tilde{\mathcal{F}}$ . Note that the right-hand side of (22) is in the desired MLE form, with ground-truth conditional density  $\Pr_{f^*}(y, s', a' | s, a)$  specified by some  $f^* \in \tilde{\mathcal{F}}$ , thanks to the realizability assumption (Assumption 3). Now, using the PAC bound for MLE shown in Agarwal et al. (2020) (see Lemma 19), we have

$$\sum_{i=1}^N \mathbb{E}_{(s_i, a_i) \sim d^{\pi_b}} \left[ \left\| \Pr_f(\cdot, \cdot | s_i, a_i) - \Pr_{f^*}(\cdot, \cdot | s_i, a_i) \right\|_1^2 \right] \leq 8 \log(|\mathcal{F}|/\delta)$$

Since all  $(s_i, a_i)$  pairs are sampled i.i.d. from the same distribution  $d^{\pi_b}$ , we shall further conclude that

$$\mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[ \left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1^2 \right] \leq \frac{8 \log(|\mathcal{F}|/\delta)}{N}. \quad (23)$$

We proceed to further relate (23) with the desired format. For this purpose, note that

$$\begin{aligned} & \left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1 \\ &= \left\| \Pr_{\hat{f}}(y = 1, \cdot, \cdot | s, a) - \Pr_{f^*}(y = 1, \cdot, \cdot | s, a) \right\|_1 + \left\| \Pr_{\hat{f}}(y = 0, \cdot, \cdot | s, a) - \Pr_{f^*}(y = 0, \cdot, \cdot | s, a) \right\|_1 \end{aligned} \quad (24a)$$

$$= 2 \left\| \frac{\Pr(\cdot, \cdot | s, a)}{1 + \hat{f}(s, a, \cdot, \cdot)} - \frac{\Pr(\cdot, \cdot | s, a)}{1 + f^*(s, a, \cdot, \cdot)} \right\|_1 \quad (24b)$$

$$= 2 \int \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} ds' da', \quad (24c)$$

where in (24a) we use the definition of  $L^1$ -norm; in (24b) we use the fact that

$$\Pr_{\hat{f}}(y, \cdot, \cdot | s, a) - \Pr_{f^*}(y, \cdot, \cdot | s, a) = (-1)^y \left( \frac{\Pr(\cdot, \cdot | s, a)}{1 + \hat{f}(s, a, \cdot, \cdot)} - \frac{\Pr(\cdot, \cdot | s, a)}{1 + f^*(s, a, \cdot, \cdot)} \right).$$

Now, plugging Lemma 7 and (21) into the integrand in (24c), we have

$$\begin{aligned} & \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} \\ &= \frac{|\mathbb{P}^\pi(s', a' | s, a)/P_{\text{neg}}(s', a') - \hat{f}(s, a, s', a')| \cdot \frac{1}{2}(\mathbb{P}^\pi(s', a' | s, a) + P_{\text{neg}}(s', a'))}{(1 + \hat{f}(s, a, s', a'))(1 + \mathbb{P}^\pi(s', a' | s, a)/P_{\text{neg}}(s', a'))} \end{aligned} \quad (25a)$$

$$= \frac{|\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a')\hat{f}(s, a, s', a')|}{2(1 + \hat{f}(s, a, s', a'))} \quad (25b)$$

$$\geq \frac{|\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a')\hat{f}(s, a, s', a')|}{2(1 + C_d)}, \quad (25c)$$

where we use the upper bound  $\hat{f}(s, a, s', a') = \hat{\mathbb{P}}^\pi(s', a' | s, a)/d^{\pi_b}(s', a') \leq C_d$  in (25c). Consequently,

$$\begin{aligned} & \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \\ &= \int |\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a')\hat{f}(s, a, s', a')| ds' da' \end{aligned} \quad (26a)$$

$$\leq 2(1 + C_d) \int \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} ds' da' \quad (26b)$$

$$= (1 + C_d) \left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1, \quad (26c)$$

where we use the definition  $\hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) = P_{\text{neg}}(\cdot, \cdot)\hat{f}(s, a, \cdot, \cdot)$  in (26a), (25) in (26b), and (24) in (26c). Finally,

$$\begin{aligned} & \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[ \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \\ &\leq \sqrt{\mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[ \left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1^2 \right]} \end{aligned} \quad (27a)$$

$$\leq \sqrt{(1 + C_d)^2 \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[ \left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1^2 \right]} \quad (27b)$$

$$\leq 2\sqrt{2}(1 + C_d) \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}, \quad (27c)$$

where we use Cauchy-Schwartz inequality in (27a), (26) in (27b), and (23) in (27c).  $\square$

## D Sample Complexity Guarantee

In this appendix, we derive the sample complexity guarantee for the proposed SPECTRALDICE algorithm, assuming a known bound on the representation learning error induced by the REPLEARN subroutine (see Claim 3). As discussed in the main text, the objective is to bound the estimation error  $\mathcal{E} := \hat{\rho}(\pi) - \rho(\pi)$ , which can be intuitively split into the following three terms that are easier to bound:

$$\mathcal{E} = \underbrace{\hat{\rho}(\pi) - \bar{\rho}(\pi)}_{\text{statistical error}} + \underbrace{\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi)}_{\text{dataset error}} + \underbrace{\rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi)}_{\text{representation error}}.$$



We point out that the statistical error results from replacing the expectation with empirical estimates, the dataset error comes from the offline dataset that samples transitions from the true transition kernel  $\mathbb{P}^\pi$  instead of the learned kernel  $\hat{\mathbb{P}}^\pi$ , and the representation error accounts for the error induced by plugging in the learned representation  $(\hat{\phi}, \hat{\mu}^\pi)$  instead of the ground truth  $(\phi^*, \mu^{\pi,*})$  into the DICE estimator.

As described in the proof sketch, for the rest of this appendix, we provide an upper bound for each of these three terms, and eventually conclude with an overall sample complexity guarantee.

**Representation Error.** We start by bounding the representation error term, which by intuition should be a direct consequence of the representation learning error shown in Claim 3.

**Lemma 9.** *Conditioned on the event that the inequality in Claim 3 holds, under Assumptions 1 to 3, we have*

$$\rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) \leq \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta).$$

*Proof.* By the well-known Simulation Lemma (see Lemma 20), we have

$$\begin{aligned} & \rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) \\ &= \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^\pi} \left[ \mathbb{E}_{s' \sim \hat{\mathbb{P}}(\cdot|s,a)} [V_{\hat{\mathbb{P}}}^\pi(s')] - \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)} [V_{\hat{\mathbb{P}}}^\pi(s')] \right] \end{aligned} \quad (28a)$$

$$= \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^\pi} \left[ \mathbb{E}_{(s',a') \sim \hat{\mathbb{P}}^\pi(\cdot, \cdot|s,a)} [Q_{\hat{\mathbb{P}}}^\pi(s', a')] - \mathbb{E}_{(s',a') \sim \mathbb{P}^\pi(\cdot, \cdot|s,a)} [Q_{\hat{\mathbb{P}}}^\pi(s', a')] \right] \quad (28b)$$

$$= \frac{\gamma}{1-\gamma} \int d_{\mathbb{P}}^\pi(s, a) ds da \int Q_{\hat{\mathbb{P}}}^\pi(s', a') \left( \hat{\mathbb{P}}^\pi(s', a'|s, a) - \mathbb{P}^\pi(s', a'|s, a) \right) ds' da' \quad (28c)$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \int C_\infty^\pi d_{\mathbb{P}}^{\pi_b}(s, a) ds da \int |\mathbb{P}^\pi(s', a'|s, a) - \hat{\mathbb{P}}^\pi(s', a'|s, a)| ds' da' \quad (28d)$$

$$= \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[ \|\mathbb{P}^\pi(s', a'|s, a) - \hat{\mathbb{P}}^\pi(s', a'|s, a)\|_1 \right] \quad (28e)$$

$$\leq \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta), \quad (28f)$$

where in (28a) we use the Simulation Lemma; in (28b) we use the relationship between value functions; in (28d) we plug in  $d_{\mathbb{P}}^\pi(s, a) \leq C_\infty^\pi d_{\mathbb{P}}^{\pi_b}(s, a)$  (Assumption 2) and the fact that  $Q_{\hat{\mathbb{P}}}^\pi(\cdot, \cdot) \leq \frac{1}{1-\gamma}$ ; in (28f) we use Claim 3.  $\square$

**Dataset Error.** The dataset error can be accounted for by a bounded difference in the objective function, which turns out to be another consequence of the representation learning error. For this purpose, we first show the following technical lemma that formalizes the above intuition.

**Lemma 10.**  $\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} F_1(\mathbf{x}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} F_2(\mathbf{x}, \mathbf{y}) \leq \max_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}} |F_1(\mathbf{x}, \mathbf{y}) - F_2(\mathbf{x}, \mathbf{y})|.$

*Proof.* Let  $\varepsilon := \max_{\mathbf{x}, \mathbf{y}} |F_1(\mathbf{x}, \mathbf{y}) - F_2(\mathbf{x}, \mathbf{y})|$ . Then we have

$$\min_{\mathbf{x}} \max_{\mathbf{y}} F_1(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x}} \left\{ \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}} \{F_1(\mathbf{x}, \mathbf{y}) - F_2(\mathbf{x}, \mathbf{y})\} \right\} \quad (29a)$$

$$\leq \min_{\mathbf{x}} \left\{ \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \varepsilon \right\} \quad (29b)$$

$$= \min_{\mathbf{x}} \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \varepsilon, \quad (29c)$$

where in (29a) we use the fact that  $\max_{\mathbf{y}} \{f(\mathbf{y}) + g(\mathbf{y})\} \leq \max_{\mathbf{y}} f(\mathbf{y}) + \max_{\mathbf{y}} g(\mathbf{y})$ .  $\square$

Now we are ready to show the following lemma regarding dataset error.

**Lemma 11.** *Conditioned on the event that the inequality in Claim 3 holds, under Assumptions 1 to 3, we have*

$$\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi) \leq \frac{C_\infty^\pi}{1-\gamma} \cdot \xi(|\mathcal{F}|, N, \delta).$$

*Proof.* For the sake of clarity, denote the optimization objectives of  $\bar{\rho}(\pi)$  and  $\rho_{\hat{\mathbb{P}}}(\pi)$  as follows:

$$\bar{\rho}(\pi) = \min_{\theta_Q} \max_{\omega_d} \bar{F}(\theta_Q, \omega_d), \quad \rho_{\hat{\mathbb{P}}}(\pi) = \min_{\theta_Q} \max_{\omega_d} \hat{F}(\theta_Q, \omega_d).$$

Then we can show that

$$|\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi)| = \left| \min_{\theta_Q} \max_{\omega_d} \bar{F}(\theta_Q, \omega_d) - \min_{\theta_Q} \max_{\omega_d} \hat{F}(\theta_Q, \omega_d) \right| \leq \left| \bar{F}(\theta_Q, \omega_d) - \hat{F}(\theta_Q, \omega_d) \right| \quad (30a)$$

$$= \left| \int d_{\mathbb{P}^b}^{\pi}(s, a) \left( \mathbb{P}^{\pi}(s', a' | s, a) - \hat{\mathbb{P}}^{\pi}(s', a' | s, a) \right) \left( \hat{\mu}^{\pi}(s, a)^{\top} \omega_d \right) \cdot \right. \\ \left. \left( r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q \right) ds da ds' da' \right| \quad (30b)$$

$$\leq \int d_{\mathbb{P}^b}^{\pi}(s, a) \left| \mathbb{P}^{\pi}(s', a' | s, a) - \hat{\mathbb{P}}^{\pi}(s', a' | s, a) \right| \cdot \left| \hat{\mu}^{\pi}(s, a)^{\top} \omega_d \right| \cdot \\ \left| r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q \right| ds da ds' da' \quad (30c)$$

$$\leq \mathbb{E}_{(s, a) \sim d_{\mathbb{P}^b}^{\pi}} \left[ \left\| \mathbb{P}^{\pi}(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a) \right\|_1 \cdot C_{\infty}^{\pi} \cdot \frac{1}{1-\gamma} \right] \quad (30d)$$

$$= \frac{C_{\infty}^{\pi}}{1-\gamma} \cdot \xi(|\mathcal{F}|, N, \delta), \quad (30e)$$

where in (30a) we use Lemma 10; in (30c) we use the integral triangle inequality; in (30d) we plug in  $|\hat{\mu}^{\pi}(s, a)^{\top} \omega_d| \leq C_{\infty}^{\pi}$  and  $|r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q| \leq \frac{1}{1-\gamma}$  (see Remark 2); in (30e) we use Claim 3.  $\square$

**Statistical Error.** Finally, the statistical error is caused by replacing the expectations with their empirical estimations, which can be bounded by Hoeffding's concentration inequality (see Lemma 16).

**Lemma 12.** *Under Assumptions 1 to 3, with probability at least  $1 - \delta$ , we have*

$$\hat{\rho}(\pi) - \bar{\rho}(\pi) \leq \frac{C_{\infty}^{\pi}}{1-\gamma} \sqrt{\frac{\log(1/2\delta)}{2N}}.$$

*Proof.* For clarity, label the samples in  $\mathcal{D}$  as  $\mathcal{D} = \{(s_i, a_i, s'_i, a'_i) \mid i \in [N]\}$ , and define

$$F(s, a, s', a') := (\hat{\mu}^{\pi}(s, a)^{\top} \omega_d) (r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q).$$

Note that  $|\hat{\mu}^{\pi}(s, a)^{\top} \omega_d| \leq C_{\infty}^{\pi}$  and  $|r(s, a) + \gamma \hat{\phi}(s', a')^{\top} \theta_Q - \hat{\phi}(s, a)^{\top} \theta_Q| \leq \frac{1}{1-\gamma}$  (see Remark 2), we have

$$|F(s, a, s', a')| \leq \frac{C_{\infty}^{\pi}}{1-\gamma}, \quad \forall s, a, s', a'.$$

Therefore, by Hoeffding's inequality (see Lemma 16), we conclude that

$$\Pr \left[ \left| \frac{1}{N} \sum_{i=1}^N F(s_i, a_i, s'_i, a'_i) - \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), \\ s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')}} [F(s, a, s', a')] \right| > t \right] \leq 2 \exp \left( -\frac{2Nt^2}{4(C_{\infty}^{\pi})^2/(1-\gamma)^2} \right).$$

Or equivalently, with probability at least  $1 - \delta$ , we have

$$\left| \hat{\mathbb{E}}_{\substack{(s, a, s') \sim \mathcal{D}, \\ a' \sim \pi(\cdot|s')}} [F(s, a, s', a')] - \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), \\ s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')}} [F(s, a, s', a')] \right| \leq \frac{C_{\infty}^{\pi}}{1-\gamma} \sqrt{\frac{\log(1/2\delta)}{2N}}.$$

Finally, the conclusion follows from Lemma 10 using the same argument as above.  $\square$

**Conclusion.** Now we are ready to prove the Main Theorem.

**Theorem 4.** *Suppose Claim 3 holds for the  $\text{REPLEARN}(\mathcal{F}, \mathcal{D}, \pi)$  subroutine. Then under Assumptions 1 to 3, with probability at least  $1 - \delta$ , we have*

$$\mathcal{E} \leq \frac{C_{\infty}^{\pi}}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}} + \frac{C_{\infty}^{\pi}}{(1-\gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta/2).$$

*Proof.* Consider the following high-probability events:

$$\mathcal{C}_1 : \mathbb{E}_{(s, a) \sim d_{\mathbb{P}^b}^{\pi}} \left[ \left\| \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a) - \mathbb{P}^{\pi}(\cdot, \cdot | s, a) \right\|_1 \right] \leq \xi(|\mathcal{F}|, N, \delta/2),$$

$$\mathcal{C}_2 : \left| \hat{\mathbb{E}}_{\substack{(s, a, s') \sim \mathcal{D}, \\ a' \sim \pi(\cdot|s')}} [F(s, a, s', a')] - \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), \\ s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')}} [F(s, a, s', a')] \right| \leq \frac{C_{\infty}^{\pi}}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}}.$$

As per Claim 3 and Lemma 12, we know  $\Pr[\mathcal{C}_i] \geq 1 - \delta/2$  ( $i = 1, 2$ ). Hence by Union Bound we have

$$\Pr[\mathcal{C}_1 \cap \mathcal{C}_2] \geq 1 - \delta.$$

On the other hand, conditioned on  $\mathcal{C}_1 \cap \mathcal{C}_2$ , Lemma 9, Lemma 11 and Lemma 12 in combination guarantee

$$\mathcal{E} = \hat{\rho}(\pi) - \bar{\rho}(\pi) + \bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi) + \rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi)$$

$$\begin{aligned}
 &\leq \frac{C_\infty^\pi}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}} + \frac{C_\infty^\pi}{1-\gamma} \cdot \xi(|\mathcal{F}|, N, \delta/2) + \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \xi(|\mathcal{F}|, N, \delta/2) \\
 &= \frac{C_\infty^\pi}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}} + \frac{C_\infty^\pi}{(1-\gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta/2)
 \end{aligned}$$

This completes the proof.  $\square$

For completeness, we also include the corollaries of the Main Theorem that characterize the overall sample complexity of our SPECTRALDICE algorithm using OLS and NCE representation learning methods.

**Corollary 13** (sample complexity of OLS-based SPECTRALDICE). *Under Assumptions 1 to 3 and the additional Assumption 4 for regularity, let  $(\hat{\phi}, \hat{\mu}^\pi)$  be the solution to the OLS problem (14). Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have*

$$\mathcal{E} \leq \frac{C_\infty^\pi}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}} + \frac{C_\infty^\pi \sqrt{C_{\mathbb{P}} C_{\text{reg}}}}{(1-\gamma)^2} \cdot \sqrt{\frac{\log(2|\mathcal{F}|/\delta)}{N}} \lesssim \frac{1}{(1-\gamma)^2} \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}},$$

where  $C_{\text{reg}} = \frac{4}{3}\sqrt{C_{\text{cov}}} + 8C_{\text{cov}}$  is a universal constant determined by the PAC bound for OLS.

**Corollary 14** (sample complexity of NCE-based SPECTRALDICE). *Under Assumptions 1 to 3 and the additional Assumption 5 for regularity, let  $(\hat{\phi}, \hat{\mu}^\pi)$  be the solution to the NCE problem (17) with  $P_{\text{neg}}(\cdot, \cdot) \equiv d^{\pi_b}(\cdot, \cdot)$ . Then, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , we have*

$$\mathcal{E} \leq \frac{C_\infty^\pi}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}} + \frac{2\sqrt{2}C_\infty^\pi(1+C_d)}{(1-\gamma)^2} \cdot \sqrt{\frac{\log(2|\mathcal{F}|/\delta)}{N}} \lesssim \frac{1}{(1-\gamma)^2} \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}.$$

*Remark 5* (Sampling the dataset). Throughout this paper, we have been slightly abusing the notation  $(s, a, s') \sim \mathcal{D}$ , which is a little subtle in practice since only trajectories (rather than transitions) are collected. To ensure the correct data distribution  $d^{\mathcal{D}}(s, a) = d^{\pi_b}(s, a)$ , we shall first randomly sample the trajectories, within which we sample each transition  $(s_t, a_t, s_{t+1}, a_{t+1})$  with probability  $(1-\gamma)\gamma^t$ .

## E Technical Lemmas

In this final appendix, we include all the technical lemmas used in the previous sections.

### E.1 $f$ -Divergence

**Definition 1** ( $f$ -divergence). Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probabilities distribution over a sample space  $\mathcal{X}$ , such that  $\mathbb{P}$  is absolutely continuous with respect to  $\mathbb{Q}$ . Given a convex function  $f : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  such that  $f(1) = 0$  and  $f(0) := \lim_{t \rightarrow 0^+} f(t)$ . Then the  $f$ -divergence of  $\mathbb{P}$  with respect to  $\mathbb{Q}$  is defined as

$$D_f(\mathbb{P} \parallel \mathbb{Q}) := \int_{\mathcal{X}} f\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{Q}.$$

The following *variational representation* of  $f$ -divergences is well-known in literature.

**Lemma 15** (variational representation using Fenchel conjugate). *Let  $\mathcal{F}$  denote the class of measurable real valued functions on  $\mathcal{X}$  that is absolutely integratable with respect to  $\mathbb{Q}$ . Then*

$$D_f(\mathbb{P} \parallel \mathbb{Q}) = \sup_{g \in \mathcal{F}} \left\{ \mathbb{E}_{x \sim \mathbb{P}}[g(x)] - \mathbb{E}_{x \sim \mathbb{Q}}[f_*(g(x))] \right\},$$

where  $f_*$  is the Fenchel conjugate of  $f$ . Further, if  $f$  is differentiable, then the optimal dual variable is given by

$$g^*(x) = f'\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) \implies D_f(\mathbb{P} \parallel \mathbb{Q}) = \mathbb{E}_{x \sim \mathbb{P}}\left[f'\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right)\right] - \mathbb{E}_{x \sim \mathbb{Q}}\left[f_*\left(f'\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right)\right)\right]$$

*Proof.* See Theorem 4.4 in Broniatowski and Keziou (2006).  $\square$

### E.2 Concentration Inequalities

**Lemma 16** (Hoeffding's inequality, Hoeffding (1994)). *Let  $X_1, X_2, \dots, X_N$  be i.i.d. random variables with mean  $\mu$  and taking values in  $[a, b]$  almost surely. Then for any  $\varepsilon > 0$  we have*

$$\Pr\left[\left|\frac{1}{N} \sum_{i=1}^N X_i - \mu\right| > \varepsilon\right] \leq 2 \exp\left(-\frac{2N\varepsilon^2}{(b-a)^2}\right).$$

In other words, with probability at least  $1 - \delta$ , we have

$$\left| \frac{1}{N} \sum_{i=1}^N X_i - \mu \right| \leq (b - a) \sqrt{\frac{\log(1/2\delta)}{2N}}.$$

**Lemma 17** (Bernstein's inequality, [Bernstein \(1924\)](#)). *Let  $X_1, X_2, \dots, X_N$  be i.i.d. random variables with mean  $\mu$ , variance  $\sigma^2$ , and bounded range  $|X_i - \mu| \leq B$  almost surely. Then with probability at least  $1 - \delta$ , we have*

$$\pm \left( \frac{1}{N} \sum_{i=1}^N X_i - \mu \right) \leq \sigma \sqrt{\frac{2 \log(1/\delta)}{N}} + \frac{B \log(1/\delta)}{3N}.$$

### E.3 Statistical Learning: PAC Bounds

In this section, we present the standard PAC bounds for OLS and MLE. Although these are both classic results, we fail to trace back to the original literature of the former, and thus provide a short proof here for completeness.

**Lemma 18** (PAC bound for OLS, fast rate). *Consider a regression problem over a finite family  $\mathcal{F} = \{f : \mathcal{X} \rightarrow [a, b]\}$  of bounded functions with data distribution  $(X, Y) \sim \mathcal{M}$ , where the objective is to solve for*

$$\arg \min_{f \in \mathcal{F}} \mathcal{L}(f), \quad \text{where } \mathcal{L}(f) := \mathbb{E}_{(X, Y) \sim \mathcal{M}} [(f(X) - Y)^2].$$

*Suppose the regression function  $f^*(x) := \mathbb{E}[Y | X = x] \in \mathcal{F}$  (realizability), and we have access to i.i.d. sample  $(x_i, y_i) \sim \mathcal{M}$ ,  $\forall i \in [N]$ . Let the Empirical Risk Minimization (ERM) estimator be*

$$\hat{f} := \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}(f), \quad \text{where } \hat{\mathcal{L}}(f) := \frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2.$$

*Then, with probability at least  $1 - \delta$ , the ERM estimator induces a regret that is at most*

$$\mathcal{L}(\hat{f}) \leq \mathcal{L}(f^*) + C_{\text{reg}} \frac{\log(|\mathcal{F}|/\delta)}{N}.$$

*Suppose further that the ground truth is deterministic such that  $y = f^*(x)$  for some  $f^* \in \mathcal{F}$ , in which case we have*

$$\mathcal{L}(\hat{f}) \leq C_{\text{reg}} \frac{\log(|\mathcal{F}|/\delta)}{N}.$$

*Here  $C_{\text{reg}} = 8(b - a)^2 + \frac{4}{3}(b - a)$  is a universal constant depending only on the range  $[a, b]$ .*

*Proof.* Define a random variable  $Z_i := (f(X_i) - Y_i)^2 - (f^*(X_i) - Y_i)^2$ , such that

$$\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [Z_i(f)] = \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - Y_i)^2 - (f^*(X_i) - Y_i)^2] \quad (31a)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} \left[ ((f(X_i) - f^*(X_i)) + (f^*(X_i) - Y_i))^2 - (f^*(X_i) - Y_i)^2 \right] \quad (31b)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))^2] + 2\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))(f^*(X_i) - Y_i)] \quad (31c)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))^2] =: \mathcal{E}(f), \quad (31d)$$

where in (31c) we use the following fact  $\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))(f^*(X_i) - Y_i)] = \mathbb{E}_{X_i} [(f(X_i) - f^*(X_i)) \cdot \mathbb{E}_{Y_i \sim \mathcal{M}(\cdot | X_i)} [f^*(X_i) - Y_i]] = 0$  that directly follows from the definition of  $f^*$ . Similarly, for any  $t \in [T]$ ,

$$\text{Var}_{(X_i, Y_i) \sim \mathcal{M}} [Z_i(f)] = \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [Z_i(f)^2] - (\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [Z_i(f)])^2 \quad (32a)$$

$$\leq \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} \left[ ((f(X_i) - Y_i)^2 - (f^*(X_i) - Y_i)^2)^2 \right] \quad (32b)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))^2 (f(X_i) + f^*(X_i) - 2Y_i)^2] \quad (32c)$$

$$\leq 4(b - a)^2 \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}} [(f(X_i) - f^*(X_i))^2] = 4(b - a)^2 \mathcal{E}(f). \quad (32d)$$

where in (32a) we simply drop the second term, and in (32c) we use the fact  $f(X_i) + f^*(X_i) - 2Y_i \in [-2(b - a), 2(b - a)]$  as  $f(X_i), f^*(X_i), Y_i \in [a, b]$ . Further, for any  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  and  $f \in \mathcal{F}$ , we have  $f(x) - y \in [-(b - a), b - a]$ , implying  $Z_i(f) \in [-(b - a), b - a]$  and  $\mathbb{E}[Z_i(f)] \in [-(b - a), (b - a)]$ . Therefore,  $|Z_i(f) - \mathbb{E}[Z_i(f)]| \leq 2(b - a)$ . Then by Bernstein's inequality ([Lemma 17](#)), we conclude that, with probability at least  $1 - \delta$ ,

$$\mathbb{E}[Z_i(f)] - \frac{1}{N} \sum_{i=1}^N Z_i(f) \leq \sqrt{\text{Var}[Z_i(f)]} \sqrt{\frac{2 \log(1/\delta)}{N}} + \frac{2(b - a) \log(1/\delta)}{3N}. \quad (33)$$

To proceed, plug (31) and (32) into (33), and we have

$$\mathcal{E}(f) - (\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) \leq 2(b - a) \sqrt{\mathcal{E}(f)} \sqrt{\frac{2 \log(1/\delta)}{N}} + \frac{2(b - a) \log(1/\delta)}{3N} \quad (34a)$$

$$\leq \left( \frac{1}{2} \mathcal{E}(f) + \frac{4(b - a)^2 \log(1/\delta)}{N} \right) + \frac{2(b - a) \log(1/\delta)}{3N}, \quad (34b)$$



where in (34a) we apply the AM-GM inequality. Finally, we rearrange the terms to obtain

$$\mathcal{E}(f) \leq 2(\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) + \frac{C_{\text{reg}} \log(1/\delta)}{N} \quad (35)$$

for any fixed  $f \in \mathcal{F}$ , with probability at least  $1 - \delta$ . Finally, we take the union bound with respect to all  $f \in \mathcal{F}$ , such that with probability at least  $1 - \delta$ , we have

$$\mathcal{E}(f) \leq 2(\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) + \frac{C_{\text{reg}} \log(|\mathcal{F}|/\delta)}{N}, \quad \forall f \in \mathcal{F}. \quad (36)$$

In particular, (36) also applies to the ERM estimator  $\hat{f}$ , which gives

$$\mathcal{E}(\hat{f}) \leq 2(\hat{\mathcal{L}}(\hat{f}) - \hat{\mathcal{L}}(f^*)) + \frac{C_{\text{reg}} \log(|\mathcal{F}|/\delta)}{N} \leq \frac{C_{\text{reg}} \log(|\mathcal{F}|/\delta)}{N}. \quad (37)$$

Here we use the inequality  $\hat{\mathcal{L}}(\hat{f}) \leq \hat{\mathcal{L}}(f^*)$ , as  $\hat{f}$  minimizes  $\hat{\mathcal{L}}(\cdot)$  within  $\mathcal{F}$ . This completes the proof.  $\square$

**Lemma 19** (PAC bound for MLE, Agarwal et al. (2020)). *Consider a conditional probability estimation problem over a finite family  $\mathcal{F} = \{f : (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}\}$ , where the objective is to estimate  $f^*(x, y) := \mathbb{P}(y|x)$ . Suppose the ground truth  $f^* \in \mathcal{F}$  (realizability), and we have access to (potentially correlated) samples  $\{(x_i, y_i) \mid i \in [N]\}$  such that  $x_i \sim \mathcal{D}_i$  ( $\mathcal{D}_i$  is allowed to depend on  $(x_{1:i-1}, y_{1:i-1})$ , forming a martingale process) and  $y_i \sim \mathbb{P}(\cdot|x_i)$ . Let the Maximum Likelihood Estimator (MLE) be*

$$\hat{f} := \arg \max_{f \in \mathcal{F}} \sum_{i=1}^N \log f(x_i, y_i).$$

Then, with probability at least  $1 - \delta$ , the error of the MLE estimator is bounded as follows:

$$\sum_{i=1}^N \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \|\hat{f}(x, \cdot) - f^*(x, \cdot)\|_1^2 \right] \leq 8 \log(|\mathcal{F}|/\delta).$$

Specifically, when  $\{(x_i, y_i) \mid i \in [N]\}$  are i.i.d. samples from a dataset  $\mathcal{D}$ , we have

$$\mathbb{E}_{x \sim \mathcal{D}} \left[ \|\hat{f}(x, \cdot) - f^*(x, \cdot)\|_1^2 \right] \leq \frac{8 \log(|\mathcal{F}|/\delta)}{N}.$$

#### E.4 Simulation Lemma in MDPs

The following Simulation Lemma is a simplified version of Lemma 21 in Uehara et al. (2021).

**Lemma 20** (Simulation Lemma). *Given two MDPs  $(\mathbb{P}, r)$  and  $(\hat{\mathbb{P}}, r)$ , for any policy  $\pi \in \Pi$ , we have*

$$\rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) = \frac{\gamma}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} \left[ \mathbb{E}_{s' \sim \hat{\mathbb{P}}(\cdot|s,a)} [V_{\hat{\mathbb{P}}}^{\pi}(s')] - \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)} [V_{\mathbb{P}}^{\pi}(s')] \right].$$

*Proof.* Note that, for any uniformly bounded function  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , we have

$$\begin{aligned} & \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [f(s, a)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{\pi, \mathbb{P}} \left[ \sum_{t=0}^{\infty} (\gamma^t f(s_t, a_t) - \gamma^{t+1} f(s_{t+1}, a_{t+1})) \mid s_0 \sim \mu_0, a_0 \sim \pi(\cdot|s_0) \right] \\ &= \frac{1}{1 - \gamma} \sum_{s,a} f(s, a) \cdot \mathbb{E}_{\pi, \mathbb{P}} \left[ \sum_{t=0}^{\infty} (\gamma^t \mathbb{1}\{s_t = s, a_t = a\} - \gamma^{t+1} \mathbb{1}\{s_{t+1} = s, a_{t+1} = a\}) \mid s_0 \sim \mu_0, a_0 \sim \pi(\cdot|s_0) \right] \\ &= \frac{1}{1 - \gamma} \sum_{s,a} f(s, a) \cdot \left( d_{\mathbb{P}}^{\pi}(s, a) - \gamma \sum_{\tilde{s}, \tilde{a}} d_{\mathbb{P}}^{\pi}(\tilde{s}, \tilde{a}) \mathbb{P}^{\pi}(s, a | \tilde{s}, \tilde{a}) \right) \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} [f(s, a) - \gamma \mathbb{E}_{(s',a') \sim \mathbb{P}^{\pi}(\cdot, \cdot | s, a)} [f(s', a')]]. \end{aligned}$$

Therefore, since  $\rho_{\mathbb{P}}(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} [r(s, a)]$  and  $\rho_{\hat{\mathbb{P}}}(\pi) = \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q_{\hat{\mathbb{P}}}^{\pi}(s, a)]$ , we have

$$\begin{aligned} \rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) &= \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q_{\hat{\mathbb{P}}}^{\pi}(s, a)] - \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} [r(s, a)] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} [Q_{\hat{\mathbb{P}}}^{\pi}(s, a) - \gamma \mathbb{E}_{(s',a') \sim \mathbb{P}^{\pi}(\cdot, \cdot | s, a)} [Q_{\hat{\mathbb{P}}}^{\pi}(s', a')] - r(s, a)] \\ &= \frac{\gamma}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\pi}^{\mathbb{P}}} \left[ \mathbb{E}_{(s',a') \sim \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a)} [Q_{\hat{\mathbb{P}}}^{\pi}(s', a')] - \mathbb{E}_{(s',a') \sim \mathbb{P}^{\pi}(\cdot, \cdot | s, a)} [Q_{\hat{\mathbb{P}}}^{\pi}(s', a')] \right], \end{aligned}$$

where in the last equality we plug in the Bellman equation  $Q_{\hat{\mathbb{P}}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{(s',a') \sim \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a)} [Q_{\hat{\mathbb{P}}}^{\pi}(s', a')]$ . Finally, we leverage the relationship between  $Q$ - and  $V$ -functions to complete the proof.  $\square$