

Primal-Dual Spectral Representation for Off-policy Evaluation

Anonymous CPAL submission

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

18 1 Introduction

19 The past decade has witnessed the ubiquitous success of reinforcement learning (RL) across vari-
20 ous domains. Despite the original rationale that RL agents should learn a reward-maximizing policy
21 from continuous interactions with the environment, there also exist a wide range of applicational
22 scenarios where *online* interaction with the environment may be expensive, inefficient, risky, uneth-
23 ical, and/or even infeasible, examples of which include robotics [1, 2], autonomous driving [3, 4],
24 healthcare [5, 6], education [7, 8], dialogue systems [9, 10] and recommendation systems [11, 12].
25 These application scenarios motivate the study of *offline* RL, where the learning agent only has access
26 to historical data collected by a separate behavior policy.

27 Off-policy evaluation (OPE) is one of the most fundamental problems in offline RL that aims at
28 estimating the expected cumulative reward of a given target policy using only historical data col-
29 lected by a different, potentially unknown behavior policy. In the past decade, various off-policy
30 performance estimators have been proposed [13–16]. However, these estimators generally suffer
31 from the *curse of horizon* [17]—step-wise variances accumulate in a multiplicative way, resulting in
32 prohibitively high trajectory variances and thus unreliable estimators. The recently proposed Dis-
33 tribution Correction Estimation (DICE) family of estimators have advanced the state of the art in
34 OPE, leveraging the primal-dual formulation of policy evaluation for a saddle-point optimization
35 approach that directly estimates the stationary distribution correction ratio, and hence breaking the
36 curse of horizon [18, 19].

37 Nevertheless, as systems scale up in terms of the size of state-action spaces, the saddle-point op-
38 timization in the formulation of DICE estimators become increasingly challenging to solve. Such
39 *curse of dimensionality* is common for RL methods in general, and people have been working to allevi-
40 ate the computational burden by exploiting function approximators. However, many known func-
41 tion approximators require additional assumptions to ensure computational and statistical prop-
42 erties [20–25], which may not be easily satisfiable in practice. Moreover, the induced optimiza-

Submitted to Second Conference on Parsimony and Learning (CPAL 2025). Do not distribute.

43 tion upon function approximators may be difficult to solve [26–28]. In particular, under a generic
44 neural network parametrization, computing the DICE estimator [29] requires solving non-convex
45 non-concave saddle-point optimizations, which is known to be NP-hard in theory and also yields
46 unstable performance in practice, and is therefore regarded as intractable.

47 This dilemma brings up a very natural question:

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

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

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

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

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

70 1.1 Related Work

71 **Off-Policy Evaluation (OPE).** Off-policy evaluation has long been an active field of RL research.
72 In the case where the behavior policy is known, various off-policy performance estimators have been
73 proposed, including direct method (DM) estimators [30, 31], importance sampling (IS) estimators
74 [13, 14], doubly-robust (DR) estimators [15, 16, 32] and other mixed-type estimators [24, 33, 34],
75 which generally suffer from the *curse of dimension*. In an effort to settle this issue, there is also abun-
76 dant literature on estimating the correction ratio of the stationary distribution [17, 35], among which
77 the distribution correction estimation (DICE) family of estimators are the state of the art that lever-
78 age a novel primal-dual formulation of OPE to eliminate the curse of horizon, and in the meantime,
79 allow unknown behavior policies [18, 19, 29, 36, 37]. However, as discussed above, the induced
80 saddle-point optimization becomes unstable with neural networks, impeding the practical applica-
81 tion of DICE estimators.

82 **Spectral Representation in MDPs.** Spectral decomposition of the transition kernel is known to
83 induce a linear structure of Q -functions, which enables the design of provably efficient algorithms
84 assuming known (primal) spectral feature maps [38–40]. These algorithms break the curse of di-
85 mensionality in the sense that their computation or sample complexity is independent of the size of
86 the state-action space, but rather, only depends polynomially on the feature space dimension, the
87 intrinsic dimension of the problem.

88 With the growing interest in spectral structures of MDPs, representation learning for RL has recently
89 attracted much theory-oriented attention in the online setting [41, 42]. Practical representation-
90 based online RL algorithms have been designed via kernel techniques [43, 44], latent variable mod-

91 els [45, 46], contrastive learning [47, 48], and diffusion score matching [49]. Recently, a unified
 92 representation learning framework is proposed from a novel viewpoint that leverages the spectral
 93 decomposition of the transition operator [40].

94 Spectral representations have also been exploited in the offline setting [50–52], where the temporal
 95 difference algorithm is applied in the linear space induced by the primal spectral feature for esti-
 96 mating Q -functions. The linear structure of the occupancy measure induced by the dual spectral
 97 feature is recently utilized in Huang et al. [53], which leads to an offline RL algorithm for station-
 98 ary density ratio estimation. Although the algorithm is theoretically sound, the stationary density
 99 ratio breaks the linearity in occupancy, and hence the algorithm is not computationally efficient. As
 100 far as we know, there is no such offline RL algorithm that efficiently utilizes both primal and dual
 101 representations.

102 2 Preliminaries

103 **Notations.** Denote by $\|\cdot\|_p$ the p -norm of vectors or the L^p -norm of functionals, and by $\langle \mathbf{x}, \mathbf{y} \rangle =$
 104 $\mathbf{x}^\top \mathbf{y}$ the Euclidean inner product of vectors \mathbf{x} and \mathbf{y} . Denote by $\widehat{\mathbb{E}}_{d^{\mathcal{D}}}[\cdot]$ the empirically approximated
 105 expectation using samples from dataset $\mathcal{D} \sim d^{\mathcal{D}}$. Denote by $\Delta(S)$ the set of distributions over set S ,
 106 the element of which shall be regarded as densities whenever feasible. Denote the indicator function
 107 by $\mathbb{1}\{\cdot\}$. Write $[n] := \{1, \dots, n\}$ for $n \in \mathbb{Z}_+$. Regard $f(n) \lesssim g(n)$ as $f(n) = O(g(n))$.

108 **Markov Decision Processes (MDPs).** We consider an *infinite-horizon* discounted Markov decision
 109 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 (pos-
 110 sibly 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
 111 function; $\mu_0 \in \Delta(\mathcal{S})$ is the initial state distribution, and $\gamma \in (0, 1)$ is the reward discount factor, so
 112 that the discounted cumulative reward can be defined as $\sum_{t=0}^{\infty} \gamma^t r_t$. We consider *stationary Marko-*
 113 *vian policies* $\Pi := \{\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})\}$ that admit an action distribution depending on the current state
 114 only. Given any policy $\pi \in \Pi$, let $\mathbb{E}_{\pi, \mathbb{P}}[\cdot]$ denote the expectation over the trajectory governed by π
 115 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*)
 116 *state-action occupancy measure* under policy π , i.e., the normalized discounted probability of visiting
 117 (s, a) in a trajectory induced by policy π , defined by

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

118 Similarly, let $d_{\mathbb{P}}^{\pi}(\cdot) \in \Delta(\mathcal{S})$ denote the *state occupancy measure* subject to the relation $d_{\mathbb{P}}^{\pi}(s, a) =$
 119 $d_{\mathbb{P}}^{\pi}(s) \pi(a|s)$. Further, define the state/state-action value functions (a.k.a. V - and Q -functions) as:

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

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

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

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

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

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

125 **The Primal-Dual Characterization of $\rho(\pi)$.** Distribution Correction Estimation (DICE) [29] is a
 126 primal-dual-based method that evaluates the value of a given target policy π in the offline setting,
 127 using the linear programming (LP) formulation of policy values [54]. Specifically, it is known that

128 we can equivalently characterize $\rho(\pi)$ defined in (1) by the *primal LP*:

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

129 Further, it can be shown that strong duality holds in (1), with Lagrangian multipliers exactly the
130 state-action occupancy measures $d^\pi(\cdot, \cdot)$. Then we characterize $\rho(\pi)$ by the following *primal-dual LP*:

$$\min_{Q(\cdot, \cdot)} \max_{d(\cdot, \cdot)} (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \mathbb{E}_{(s, a) \sim d^\pi(\cdot, \cdot)} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right]. \quad (2)$$

131 We highlight that this primal-dual LP formulation is favored in the offline RL setting in that histori-
132 cal experiences can be utilized to empirically approximate the expectations in (2) after some simple
133 change-of-variables. In particular, for any measurable function $f(s, a)$, the importance sampling
134 (IS) estimator for the expected value of $f(s, a)$ is given by

$$\mathbb{E}_{(s, a) \sim d^\pi} [f(s, a)] = \mathbb{E}_{(s, a) \sim d^{\pi_b}} \left[\frac{d^\pi(s, a)}{d^{\mathcal{D}}(s, a)} \cdot f(s, a) \right], \quad (3)$$

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

136 The DICE family estimators [18, 55, 56] is designed by plugging the IS expectation estimator (3) into
137 (2), such that the stationary distribution correction ratio $\zeta(\cdot, \cdot)$ is parameterized along with the Q -
138 function to formulate an optimization, with various regularization available [57]. It is evident that
139 the DICE family estimators are applicable to the offline RL setting with unknown behavior policy.

140 **Spectral Representation.** We can always perform spectral decomposition of the dynamic operator
141 to obtain a spectral representation of *any* MDP [40]. In particular, *low-rank MDPs* refer to such MDPs
142 with intrinsic finite-rank spectral representation structures that enable scalable RL algorithms, and
143 are thus of theoretical interest [38, 58]. Formally, \mathcal{M} is said to be a *low-rank MDP* if there exists a
144 *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 $\mathbb{P}(s'|s, a) =$
145 $\langle \phi(s, a), \tilde{\mu}(s') \rangle, r(s, a) = \langle \phi(s, a), \theta_r \rangle$, for any $s, s' \in \mathcal{S}, a \in \mathcal{A}$. Here both the primal feature ϕ and
146 the dual features $\tilde{\mu}, \theta_r$ are assumed to be unknown, and thus must be learned from data [41, 42].

147 Unfortunately, it is revealed in Ren et al. [40], Zhang et al. [48] that learning the features of a low-
148 rank MDP is difficult from the unnormalized density fitting point of view. To settle this tractability
149 issue, the above papers propose a reparameterization of the dual feature as $\tilde{\mu}(\cdot) = q(\cdot)\mu(\cdot)$, where
150 we introduce an auxiliary distribution $q(\cdot) \in \Delta(\mathcal{S})$ that will be specified later. Therefore, we will
151 stick to the following spectral decomposition of the transition kernel in this paper:

$$\mathbb{P}(s'|s, a) = \langle \phi(s, a), q(s')\mu(s') \rangle, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}. \quad (4)$$

152 Under such reparameterization, it is known that the spectral representaton can be learned efficiently.

153 Additionally, we also assume μ_0 to be linearly representable in the dual feature space.

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

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

158 For satisfactory performance, it is important that the behavior policy provides sufficient data cov-
159 erage for the frequent transitions experienced by policy π . Specifically, we assume the occupancy
160 ratio between π and π_b satisfies the following regularity assumption.

161 **Assumption 2** (concentratability). $\frac{d^\pi(s, a)}{d^{\pi_b}(s, a)} \leq C_\infty^\pi, \forall s \in \mathcal{S}, a \in \mathcal{A}$.

162 We point out that the concentratability assumption is standard in offline RL literature [22, 59], and
163 is also implicitly enforced in recent work like Huang et al. [53] (see Definition 1 therein). We are

164 aware that the coefficient C_∞^π can potentially be translated into different feature-related constants
 165 [42], which does not change the asymptotics of sample complexity, yet only adds to the technical
 166 complexity. For clarity, we will stick to the simple Assumption 2 in this paper.

167 3 SPECTRALDICE: OPE using Primal-Dual Spectral Representation

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

175 3.1 Primal-Dual Spectral Representation

176 We start by specifying the primal-dual spectral representation used in SPECTRALDICE. At first glance,
 177 it may seem natural to directly learn the spectral representation of \mathbb{P} as defined in (4). However,
 178 it turns out that this naive approach includes the target policy π in the linear representation of
 179 $d^\pi(\cdot, \cdot)$, which in turn induces a complicated representation for the stationary distribution correc-
 180 tion ratio $\zeta(\cdot, \cdot)$ [53], and thus, leads to an intractable optimization (2) for the computation of the DICE
 181 estimator.

182 The above challenge inspires us to properly reparameterize the spectral decomposition (4). Specif-
 183 ically, since we only work with a fixed target policy π for off-policy evaluation, we shall con-
 184 sider the following alternative representation of the state-action transition kernel $\mathbb{P}^\pi(s', a'|s, a) :=$
 185 $\mathbb{P}(s'|s, a)\pi(a'|s')$:

$$\mathbb{P}^\pi(s', a'|s, a) = \left\langle \phi(s, a), q(s')\pi_b(a'|s') \underbrace{\frac{\pi(a'|s')}{\pi_b(a'|s')}}_{\boldsymbol{\mu}^\pi(s', a')} \boldsymbol{\mu}(s') \right\rangle. \quad (5)$$

186 Note that Assumption 2 guarantees a non-zero denominator when the nominator is non-zero. We
 187 refer to (5) as the *primal-dual spectral representation* of the state-action) transition kernel \mathbb{P}^π , where
 188 $\phi(\cdot, \cdot)$ and $\boldsymbol{\mu}^\pi(\cdot, \cdot)$ are still called *primal* and *dual* spectral features, respectively. The superscript π of
 189 the dual spectral feature emphasizes its dependence on the target policy.

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

194 **Lemma 1.** *With primal-dual spectral representation (5), the Q -function $Q^\pi(\cdot, \cdot)$ is linearly representable in*
 195 *the primal feature space with cofactor $\boldsymbol{\theta}_Q^\pi \in \mathbb{R}^d$:*

$$Q^\pi(s, a) = \langle \phi(s, a), \boldsymbol{\theta}_Q^\pi \rangle, \quad \forall s \in S, a \in \mathcal{A}. \quad (6)$$

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

$$d^\pi(s, a) = q(s)\pi_b(a|s) \langle \boldsymbol{\mu}^\pi(s, a), \boldsymbol{\omega}_d^\pi \rangle, \quad \forall s \in S, a \in \mathcal{A}.$$

198 *Specifically, when the auxiliary distribution $q(\cdot)$ is selected as the state-occupancy measure $d^{\pi_b}(\cdot)$ of the be-*
 199 *havior policy π_b , the stationary distribution correction ratio can also be linearly represented as:*

$$\zeta(s, a) = \frac{d^\pi(s, a)}{q(s)\pi_b(a|s)} = \langle \boldsymbol{\mu}^\pi(s, a), \boldsymbol{\omega}_d^\pi \rangle. \quad (7)$$

200 *Proof.* Note that the original dual feature in (4) can be restored by $\boldsymbol{\mu}(s') = \frac{\pi_b(a'|s')}{\pi(a'|s')} \boldsymbol{\mu}^\pi(s', a')$ for any
 201 $a' \in \mathcal{A}$. Then by Bellman recursive equation we have:

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

202 Similarly, by the recursive property of d^π we have:

$$\begin{aligned} d^\pi(s, a) &= (1 - \gamma) \mu_0(s) \pi(a|s) + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \mathbb{P}^\pi(s, a | \tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a} \\ &= (1 - \gamma) q(s) \langle \pi_b(a|s) \boldsymbol{\mu}^\pi(s, a), \boldsymbol{\omega}_0 \rangle + \gamma \left\langle q(s) \pi_b(a|s) \boldsymbol{\mu}^\pi(s, a), \int d^\pi(\tilde{s}, \tilde{a}) \boldsymbol{\phi}(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a} \right\rangle \\ &= \left\langle q(s) \pi_b(a|s) \boldsymbol{\mu}^\pi(s, a), \underbrace{(1 - \gamma) \boldsymbol{\omega}_0 + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \boldsymbol{\phi}(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a}}_{\boldsymbol{\omega}_d^\pi} \right\rangle, \end{aligned}$$

203 where we use the initial representation (Assumption 1) and the fact that $\pi(a|s) \boldsymbol{\mu}(s) =$
 204 $\pi_b(a|s) \boldsymbol{\mu}^\pi(s, a)$. The representation of $\zeta(\cdot, \cdot)$ is hence a direct corollary since $q(s) \pi_b(a|s) = d^{\pi_b}(s, a)$
 205 when $q(\cdot) = d^{\pi_b}(\cdot)$. \square

206 Then, using the linear spectral representations of Q and ζ in (6) and (7), we shall equivalently
 207 formulate the DICE estimator as follows.

208 **Corollary 2.** *With primal-dual spectral representation (5) where $q(\cdot) \equiv d^{\pi_b}(\cdot)$, under Assumption 1,*

$$\begin{aligned} \rho_{\mathbb{P}}(\pi) &= \min_{\boldsymbol{\theta}_Q} \max_{\boldsymbol{\omega}_d} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [\boldsymbol{\phi}(s, a)^\top \boldsymbol{\theta}_Q] \right. \\ &\quad \left. + \mathbb{E}_{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[(\boldsymbol{\mu}^\pi(s, a)^\top \boldsymbol{\omega}_d) (r(s, a) + \gamma \boldsymbol{\phi}(s', a')^\top \boldsymbol{\theta}_Q - \boldsymbol{\phi}(s, a)^\top \boldsymbol{\theta}_Q) \right] \right\}. \end{aligned} \quad (8)$$

209 The proof of Corollary 2 is deferred to Appendix B.1 due to limited space. We highlight that our
 210 new DICE formulation (8) bears several benefits:

- 211 • **Offline data compatible.** The estimator is favorable for OPE since the expectation over the (s, a, s')
 212 transition pair can be effectively approximated by samples from the offline dataset \mathcal{D} , as long as
 213 the auxiliary distribution $q(\cdot)$ is selected as the state occupancy measure d^{π_b} of the behavior policy
 214 π_b such that $\Pr[(s, a, s') \in \mathcal{D}] = q(s) \pi_b(a|s) \mathbb{P}(s'|s, a)$.
- 215 • **Optimization tractable.** Given (learned) $\boldsymbol{\phi}(s, a)$ and $\boldsymbol{\mu}^\pi(s, a)$, the saddle-point optimization in
 216 (8) is convex-concave with respect to both $\boldsymbol{\theta}_Q$ and $\boldsymbol{\omega}_d$, which perfectly bypasses the optimiza-
 217 tion difficulty in vanilla DICE estimators with neural-network-parameterized $Q^\pi(\cdot, \cdot)$ and $\zeta(\cdot, \cdot)$.
 218 Meanwhile, compared to the counterpart obtained by directly applying the naive spectral rep-
 219 resentation (4) (details of which can be found in Appendix B.2), the proposed estimator (8) is
 220 tractable in that it is free of the policy ratio $\frac{\pi(a|s)}{\pi_b(a|s)}$ that is unknown.

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

222 3.2 Spectral Representation Learning

223 In the last section, we have elaborated on how to perform OPE using off-policy data given a primal-
 224 dual spectral representation. Now it only suffices to specify how to learn such a representation,
 225 which we regard as an abstract subroutine $(\hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\mu}}^\pi) \leftarrow \text{REPLEARN}(\mathcal{F}, \mathcal{D}, \pi)$. Here \mathcal{F} denotes the col-
 226 lection of candidate representations. We highlight that our algorithm works with any representa-
 227 tion learning method that has a bounded learning error, without any further requirements on the

Algorithm 1 SPECTRALDICE: DIstribution Correction Estimation with Spectral Representation

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

1: Learn a spectral representation $(\hat{\phi}, \hat{\mu}^\pi) \leftarrow \mathbf{REPLEARN}(\mathcal{F}, \mathcal{D}, \pi)$.

2: Plug in the spectral representation $(\hat{\phi}, \hat{\mu}^\pi)$ to compute the following DICE estimator:

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

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

228 learning mechanism. Given a range of spectral representation learning methods available in liter-
 229 ature [40, 45, 48, 49], for the sake of clarity we only consider a few candidates here, while other
 230 methods may also be applicable:

231 1. **Ordinary Least Squares (OLS)**. Inspired by Ren et al. [40], an OLS objective can be constructed
 232 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
 233 transitions under behavior policy π_b , based on which we plug in (5) to obtain

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

234 which further induces the following OLS objective:

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

235 Therefore, $(\hat{\phi}, \hat{\mu}^\pi)$ can be learned by solving [40, 60]:

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

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

239 2. **Noise-Contrastive Estimation (NCE)**. NCE is a widely used method for contrastive representa-
 240 tion learning in RL [47, 48]. To learn $(\hat{\phi}, \hat{\mu}^\pi)$, we consider a binary contrastive learning objective
 241 [47]:

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

242 where P_{neg} is a negative sampling distribution.

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

245 3.3 SPECTRALDICE

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

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

251 1. **Representation learning**. We may choose any representation learning method that comes with
 252 a bounded learning error as the REPLEARN subroutine, and the overall sample complexity will
 253 depend on this choice (see Section 4).

254 **2. DICE-based policy evaluation.** With the learned representation $(\hat{\phi}, \hat{\mu}^\pi)$, we use the primal-dual
 255 DICE estimator (9) to estimate the value of the target policy π . Note that the data distribution
 256 $d^{\mathcal{D}}(s, a, s') = d^{\pi_b}(s)\pi_b(a|s)\mathbb{P}(s'|s, a)$ is exactly compatible with the formulation in (8).

257 *Remark 2* (Numerical considerations). It is known that directly solving (9) leads to potential numerical
 258 instability issues due to the objective's linearity in θ_Q and ω_d [19]. Fortunately, it is shown in
 259 Yang et al. [57] that certain regularization leads to strictly concave inner maximization while keeping
 260 the optimal *solution* unbiased (see Appendix B.3 for details). In our implementation, we append
 261 a regularizer $-\lambda \mathbb{E}_{(s,a) \sim \mathcal{D}} [f(\hat{\mu}^\pi(s, a)^\top \omega_d)]$ to the objective in (9), where f is a differentiable function
 262 with closed and convex Fenchel conjugate f_* (see Appendix E.1), and λ is a tunable constant. Fur-
 263 thermore, since $\mu^\pi(s, a)^\top \omega_d = \zeta(s, a) \leq C_\infty^\pi$ and $\phi(s, a)^\top \theta_Q = Q(s, a) \leq \frac{1}{1-\gamma}$, we also restrict θ_Q
 264 and ω_d in regions $\Theta(\hat{\phi}) = \{\theta_Q \mid 0 \leq \hat{\phi}(s, a)^\top \theta_Q \leq \frac{1}{1-\gamma}\}$ and $\Omega(\hat{\mu}^\pi) = \{\omega_d \mid \hat{\mu}^\pi(s, a)^\top \omega_d \leq C_\infty^\pi\}$,
 265 respectively.

266 4 Theoretical Guarantee

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

271 **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
 272 a valid state-action transition kernel for any $(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}$, and the ground-truth representation
 273 $(\phi^*, \mu^{\pi,*}) \in \mathcal{F}$.

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

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

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

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

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

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

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

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

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

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

288 where we introduce an auxiliary problem:

$$\bar{\rho}(\pi) = \min_{\theta_Q} \max_{\omega_d} \left\{ (1-\gamma) \mathbb{E}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} \left[\hat{\phi}(s, a)^\top \theta_Q \right] \right\}$$

$$+ \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), \\ (s', a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)}} \left[(\hat{\boldsymbol{\mu}}^\pi(s, a)^\top \boldsymbol{\omega}_d) \cdot (r(s, a) + \gamma \hat{\phi}(s', a')^\top \boldsymbol{\theta}_Q - \hat{\phi}(s, a)^\top \boldsymbol{\theta}_Q) \right].$$

289 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
 290 expectation is still taken over $(s', a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)$ rather than $\hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) = \langle \hat{\phi}(s, a), \hat{\boldsymbol{\mu}}^\pi(\cdot, \cdot) \rangle$.

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

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

293 On the other hand, the first term is only caused by replacing the expectations with their empirical
 294 estimators, which can be bounded by concentration inequalities as:

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

295 Plugging these terms back completes the proof. \square

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

299 5 Experiments

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

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

309 5.1 Continuous Environments

310 **Setting.** We start by comparing SPECTRALDICE with various baseline OPE methods in literature,
 311 including BESTDICE [57], Fitted Q Evaluation (FQE) [61], Model-Based (MB) method [62], Impor-
 312 tance Sampling (IS) method [13] and Doubly-Robust (DR) method [32]. We follow the experiment
 313 protocol in Yang et al. [57] to evaluate and compare the OPE performances of these algorithms in
 314 three continuous MuJoCo environments, namely Cartpole, Reacher and Half-Cheetah, in an in-
 315 creasing order of difficulty. In our implementation, for representation learning, we parameterize
 316 each of $\hat{\phi}$ and $\hat{\boldsymbol{\mu}}^\pi$ with a 2-layer feed-forward neural network. For the OPE step, regularizer is ap-
 317 pended to (9), and the estimated policy value is retrieved by $\hat{\rho}(\pi) = \mathbb{E}_{(s,a) \sim d^{\mathcal{D}}} [\hat{\boldsymbol{\mu}}^\pi(s, a)^\top \boldsymbol{\omega}_d \cdot r(s, a)]$
 318 (see Remark 2). Both steps are regarded as stochastic optimization problems, and are solved by
 319 stochastic gradient descent and stochastic gradient descent-ascent, respectively. Optimization hy-
 320 perparameters are selected via grid search. Performance is quantified by OPE error $|\hat{\rho}(\pi) - \rho(\pi)|$.

321 **Results.** The OPE performances of different methods in three environments are shown in Figure 1.
 322 It is observed that SPECTRALDICE achieves comparable performance in fewer optimization steps as
 323 compared to all the other baselines, and further, outperforms them in terms of both convergence rate
 324 and final estimation error in the most challenging Half-Cheetah environment. Further, although

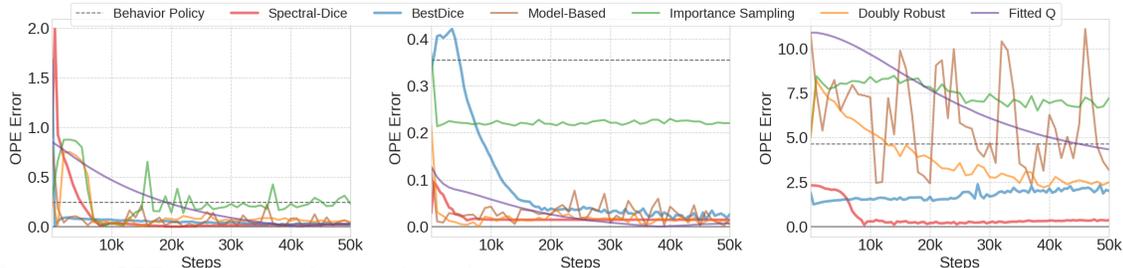


Figure 1: OPE error over the number of training steps in *Cartpole*, *Reacher* and *Half-Cheetah* environments (from left to right). Due to the use of convex-concave formulation, we can see that **SPECTRALDICE** converges faster and more stably to the target policy with a smaller OPE error in all three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

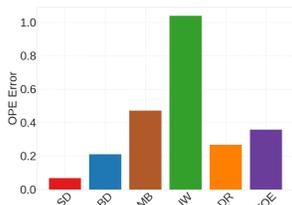


Figure 2: Averaged relative OPE error over three environments.

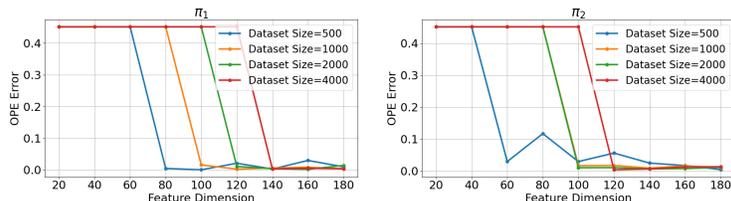


Figure 3: OPE error of **SPECTRALDICE** in *Four Rooms* with varying behavior policies (“far-away” policy π_1 vs. “similar” policy π_2), dataset sizes and feature dimensions.

325 FQE achieves an error close to **SPECTRALDICE** in simpler environments, its performance significantly
 326 degrades when the transition dynamics becomes more complex, demonstrating the importance and
 327 power of spectral representation.

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

333 For a clearer comparison, we further present the averaged relative OPE error across these three en-
 334 vironments 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 normal-
 335 ized by the value difference between the target and behavior policies. Under this metric, it becomes
 336 more evident that our method outperforms all the baselines in terms of estimation accuracy by a
 337 large margin.

338 5.2 Discrete Environment

339 **Setting.** We proceed to test our method in *Four Rooms* [63], a classical discrete environment fea-
 340 turing convenient visualization, to study the algorithm’s sensitivity for hyperparameters and il-
 341 lustrate the efficacy of representation learning. For representation learning in this tabular MDP,
 342 we perform singular value decomposition (SVD) of the matrix $\left[\frac{\mathbb{P}^\pi(s', a'|s, a)}{d^\pi(s', a')}\right]$ (indexed by (s, a) and
 343 (s', a')) and select the top d singular vectors as $\hat{\phi}(s, a)$ and $\hat{\mu}^\pi(s', a')$.

344 **Sensitivity Study.** We study the algorithm’s sensitivity with respect to behavior policy π_b , dataset
 345 size N and spectral feature dimension d by examining their impact on the OPE performance. For π_b ,
 346 we vary between two behavior policies π_1 and π_2 , where π_1 has a larger ℓ_1 -distance from the target
 347 policy than π_2 . The results are shown in Figure 3. It can be observed that the proposed algorithm
 348 is always able to achieve low OPE errors with sufficiently large feature dimensions, showcasing its
 349 wide applicability under different behavior policies, data availability and hyperparameters.

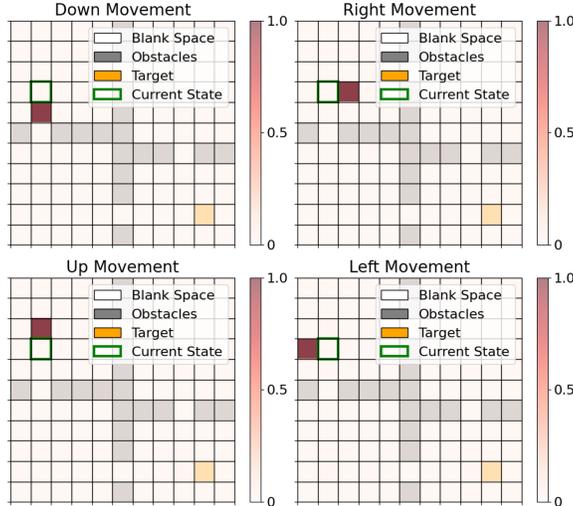


Figure 4: Visualization of the learned transition kernel for a fixed state and all the four actions.

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

355 More experimental details are deferred to Appendix A.

356 6 Conclusion

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

367 References

- 368 [1] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang,
 369 Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep
 370 reinforcement learning for vision-based robotic manipulation. In *Conference on robot learning*,
 371 pages 651–673. PMLR, 2018.
- 372 [2] Gregory Kahn, Adam Villaflor, Pieter Abbeel, and Sergey Levine. Composable action-
 373 conditioned predictors: Flexible off-policy learning for robot navigation. In *Conference on robot*
 374 *learning*, pages 806–816. PMLR, 2018.
- 375 [3] Tianyu Shi, Dong Chen, Kaian Chen, and Zhaojian Li. Offline reinforcement learning for au-
 376 tonomous driving with safety and exploration enhancement. *arXiv preprint arXiv:2110.07067*,
 377 2021.

- 378 [4] Xing Fang, Qichao Zhang, Yinfeng Gao, and Dongbin Zhao. Offline reinforcement learning for
379 autonomous driving with real world driving data. In *2022 IEEE 25th International Conference*
380 *on Intelligent Transportation Systems (ITSC)*, pages 3417–3422. IEEE, 2022.
- 381 [5] Abhyuday Jagannatha, Philip Thomas, and Hong Yu. Towards high confidence off-policy re-
382 inforcement learning for clinical applications. In *CausalML Workshop, ICML*, 2018.
- 383 [6] Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Sriniva-
384 san, Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, et al. Evaluating reinforcement
385 learning algorithms in observational health settings. *arXiv preprint arXiv:1805.12298*, 2018.
- 386 [7] Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy
387 evaluation across representations with applications to educational games. In *AAMAS*, volume
388 1077, 2014.
- 389 [8] Ahmad Slim, Husain Al Yusuf, Nadine Abbas, Chaouki T Abdallah, Gregory L Heileman, and
390 Ameer Slim. A Markov decision processes modeling for curricular analytics. In *2021 20th IEEE*
391 *international conference on machine learning and applications (ICMLA)*, pages 415–421. IEEE, 2021.
- 392 [9] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza,
393 Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement learn-
394 ing of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 395 [10] Haoming Jiang, Bo Dai, Mengjiao Yang, Tuo Zhao, and Wei Wei. Towards automatic
396 evaluation of dialog systems: A model-free off-policy evaluation approach. *arXiv preprint*
397 *arXiv:2102.10242*, 2021.
- 398 [11] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of
399 contextual-bandit-based news article recommendation algorithms. In *Proceedings of the fourth*
400 *ACM international conference on Web search and data mining*, pages 297–306, 2011.
- 401 [12] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. Top-*k*
402 off-policy correction for a reinforce recommender system. In *Proceedings of the Twelfth ACM*
403 *International Conference on Web Search and Data Mining*, pages 456–464, 2019.
- 404 [13] Josiah Hanna, Scott Niekum, and Peter Stone. Importance sampling policy evaluation with
405 an estimated behavior policy. In *International Conference on Machine Learning*, pages 2605–2613.
406 PMLR, 2019.
- 407 [14] Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Towards optimal off-policy evaluation for re-
408 inforcement learning with marginalized importance sampling. *Advances in neural information*
409 *processing systems*, 32, 2019.
- 410 [15] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning.
411 In *International conference on machine learning*, pages 652–661. PMLR, 2016.
- 412 [16] Dylan J Foster, Akshay Krishnamurthy, David Simchi-Levi, and Yunzong Xu. Offline rein-
413 forcement learning: Fundamental barriers for value function approximation. *arXiv preprint*
414 *arXiv:2111.10919*, 2021.
- 415 [17] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon:
416 Infinite-horizon off-policy estimation. *Advances in neural information processing systems*, 31, 2018.
- 417 [18] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. DualDICE: Behavior-agnostic estima-
418 tion of discounted stationary distribution corrections. *Advances in neural information processing*
419 *systems*, 32, 2019.
- 420 [19] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Al-
421 gaeDICE: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.

- 422 [20] Geoffrey J Gordon. Stable function approximation in dynamic programming. In *Machine learn-*
423 *ing proceedings 1995*, pages 261–268. Elsevier, 1995.
- 424 [21] Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire.
425 Contextual decision processes with low bellman rank are PAC-learnable. In *International Con-*
426 *ference on Machine Learning*, pages 1704–1713. PMLR, 2017.
- 427 [22] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement
428 learning. In *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- 429 [23] Wenhao Zhan, Baihe Huang, Audrey Huang, Nan Jiang, and Jason Lee. Offline reinforcement
430 learning with realizability and single-policy concentrability. In *Conference on Learning Theory*,
431 pages 2730–2775. PMLR, 2022.
- 432 [24] Pulkit Katdare, Nan Jiang, and Katherine Rose Driggs-Campbell. Marginalized importance
433 sampling for off-environment policy evaluation. In *Conference on Robot Learning*, pages 3778–
434 3788. PMLR, 2023.
- 435 [25] Fengdi Che, Chenjun Xiao, Jincheng Mei, Bo Dai, Ramki Gummadi, Oscar A Ramirez,
436 Christopher K Harris, A Rupam Mahmood, and Dale Schuurmans. Target networks and
437 over-parameterization stabilize off-policy bootstrapping with function approximation. *arXiv*
438 *preprint arXiv:2405.21043*, 2024.
- 439 [26] Justin Boyan and Andrew Moore. Generalization in reinforcement learning: Safely approxi-
440 mating the value function. *Advances in neural information processing systems*, 7, 1994.
- 441 [27] Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In
442 *Machine learning proceedings 1995*, pages 30–37. Elsevier, 1995.
- 443 [28] John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function
444 approximation. *Advances in neural information processing systems*, 9, 1996.
- 445 [29] Ofir Nachum and Bo Dai. Reinforcement learning via Fenchel-Rockafellar duality. *arXiv*
446 *preprint arXiv:2001.01866*, 2020.
- 447 [30] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with
448 Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine*
449 *Learning*, 71:89–129, 2008.
- 450 [31] Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In
451 *International Conference on Machine Learning*, pages 3703–3712. PMLR, 2019.
- 452 [32] Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning.
453 *arXiv preprint arXiv:1103.4601*, 2011.
- 454 [33] Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforce-
455 ment learning. In *International Conference on Machine Learning*, pages 2139–2148. PMLR, 2016.
- 456 [34] Nathan Kallus and Masatoshi Uehara. Double reinforcement learning for efficient off-policy
457 evaluation in markov decision processes. *Journal of Machine Learning Research*, 21(167):1–63,
458 2020.
- 459 [35] Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and Q -function learning for
460 off-policy evaluation. In *International Conference on Machine Learning*, pages 9659–9668. PMLR,
461 2020.
- 462 [36] Mengjiao Yang, Bo Dai, Ofir Nachum, George Tucker, and Dale Schuurmans. Offline policy
463 selection under uncertainty. In *International Conference on Artificial Intelligence and Statistics*,
464 pages 4376–4396. PMLR, 2022.

- 465 [37] Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estima-
466 tion of stationary values. *arXiv preprint arXiv:2002.09072*, 2020.
- 467 [38] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement
468 learning with linear function approximation. In *Conference on learning theory*, pages 2137–2143.
469 PMLR, 2020.
- 470 [39] Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels,
471 and regret bound. In *International Conference on Machine Learning*, pages 10746–10756. PMLR,
472 2020.
- 473 [40] Tongzheng Ren, Tianjun Zhang, Lisa Lee, Joseph E Gonzalez, Dale Schuurmans, and
474 Bo Dai. Spectral decomposition representation for reinforcement learning. *arXiv preprint*
475 *arXiv:2208.09515*, 2022.
- 476 [41] Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. FLAMBE: Structural
477 complexity and representation learning of low rank MDPs. *Advances in neural information pro-*
478 *cessing systems*, 33:20095–20107, 2020.
- 479 [42] Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and
480 offline RL in low-rank MDPs. *arXiv preprint arXiv:2110.04652*, 2021.
- 481 [43] Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise:
482 Provable and practical exploration for representation learning. In *Uncertainty in Artificial In-*
483 *telligence*, pages 1686–1696. PMLR, 2022.
- 484 [44] Tongzheng Ren, Zhaolin Ren, Na Li, and Bo Dai. Stochastic nonlinear control via finite-
485 dimensional spectral dynamic embedding. *arXiv preprint arXiv:2304.03907*, 2023.
- 486 [45] Tongzheng Ren, Chenjun Xiao, Tianjun Zhang, Na Li, Zhaoran Wang, Sujay Sanghavi, Dale
487 Schuurmans, and Bo Dai. Latent variable representation for reinforcement learning. *arXiv*
488 *preprint arXiv:2212.08765*, 2022.
- 489 [46] Hongming Zhang, Tongzheng Ren, Chenjun Xiao, Dale Schuurmans, and Bo Dai. Provable
490 representation with efficient planning for partially observable reinforcement learning. *arXiv*
491 *preprint arXiv:2311.12244*, 2023.
- 492 [47] Shuang Qiu, Lingxiao Wang, Chenjia Bai, Zhuoran Yang, and Zhaoran Wang. Contrastive
493 UCB: Provably efficient contrastive self-supervised learning in online reinforcement learning.
494 In *International Conference on Machine Learning*, pages 18168–18210. PMLR, 2022.
- 495 [48] Tianjun Zhang, Tongzheng Ren, Mengjiao Yang, Joseph Gonzalez, Dale Schuurmans, and
496 Bo Dai. Making linear MDPs practical via contrastive representation learning. In *International*
497 *Conference on Machine Learning*, pages 26447–26466. PMLR, 2022.
- 498 [49] Dmitry Shribak, Chen-Xiao Gao, Yitong Li, Chenjun Xiao, and Bo Dai. Diffusion spectral
499 representation for reinforcement learning. *arXiv preprint arXiv:2406.16121*, 2024.
- 500 [50] Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under
501 partial coverage. *arXiv preprint arXiv:2107.06226*, 2021.
- 502 [51] Chengzhuo Ni, Anru R Zhang, Yaqi Duan, and Mengdi Wang. Learning good state and action
503 representations via tensor decomposition. In *2021 IEEE International Symposium on Information*
504 *Theory (ISIT)*, pages 1682–1687. IEEE, 2021.
- 505 [52] Jonathan Chang, Kaiwen Wang, Nathan Kallus, and Wen Sun. Learning Bellman complete
506 representations for offline policy evaluation. In *International Conference on Machine Learning*,
507 pages 2938–2971. PMLR, 2022.

- 508 [53] Audrey Huang, Jinglin Chen, and Nan Jiang. Reinforcement learning in low-rank MDPs with
509 density features. In *International Conference on Machine Learning*, pages 13710–13752. PMLR,
510 2023.
- 511 [54] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
512 Wiley & Sons, 2014.
- 513 [55] Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Alekh Agarwal, and Wen
514 Sun. Efficient reinforcement learning in block MDPs: A model-free representation learning
515 approach. In *International Conference on Machine Learning*, pages 26517–26547. PMLR, 2022.
- 516 [56] Bo Dai, Ofir Nachum, Yinlam Chow, Lihong Li, Csaba Szepesvári, and Dale Schuurmans.
517 CoinDICE: Off-policy confidence interval estimation. *Advances in neural information process-*
518 *ing systems*, 33:9398–9411, 2020.
- 519 [57] Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy evaluation
520 via the regularized Lagrangian. *Advances in Neural Information Processing Systems*, 33:6551–
521 6561, 2020.
- 522 [58] Hengshuai Yao, Csaba Szepesvári, Bernardo Avila Pires, and Xinhua Zhang. Pseudo-MDPs
523 and factored linear action models. In *2014 IEEE Symposium on Adaptive Dynamic Programming*
524 *and Reinforcement Learning (ADPRL)*, pages 1–9. IEEE, 2014.
- 525 [59] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of*
526 *Machine Learning Research*, 9(5), 2008.
- 527 [60] Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-
528 supervised deep learning with spectral contrastive loss. *Advances in Neural Information Process-*
529 *ing Systems*, 34:5000–5011, 2021.
- 530 [61] Ilya Kostrikov and Ofir Nachum. Statistical bootstrapping for uncertainty estimation in off-
531 policy evaluation. *arXiv preprint arXiv:2007.13609*, 2020.
- 532 [62] Michael R Zhang, Tom Le Paine, Ofir Nachum, Cosmin Paduraru, George Tucker, Ziyu Wang,
533 and Mohammad Norouzi. Autoregressive dynamics models for offline policy evaluation and
534 optimization. *arXiv preprint arXiv:2104.13877*, 2021.
- 535 [63] Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A
536 framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):
537 181–211, 1999.
- 538 [64] Justin Fu, Mohammad Norouzi, Ofir Nachum, George Tucker, Ziyu Wang, Alexander Novikov,
539 Mengjiao Yang, Michael R Zhang, Yutian Chen, Aviral Kumar, et al. Benchmarks for deep off-
540 policy evaluation. *arXiv preprint arXiv:2103.16596*, 2021.
- 541 [65] Michel Broniatowski and Amor Keziou. Minimization of ϕ -divergences on sets of signed mea-
542 sures. *Studia Scientiarum Mathematicarum Hungarica*, 43(4):403–442, 2006.
- 543 [66] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The col-*
544 *lected works of Wassily Hoeffding*, pages 409–426, 1994.
- 545 [67] Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of
546 laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.

APPENDIX

548 A More Experimental Results

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

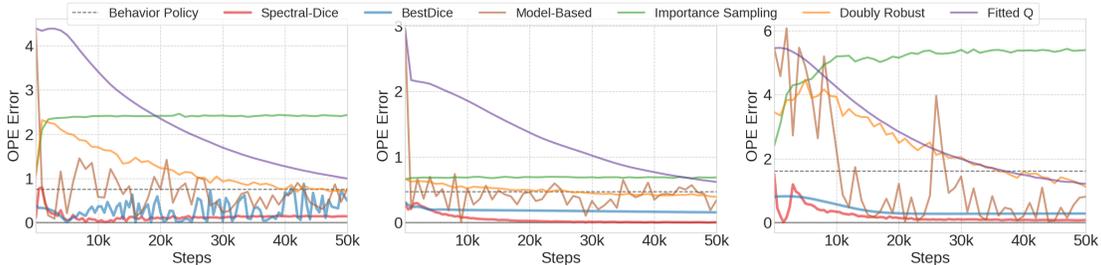


Figure 5: OPE error over the number of training steps in Walker2d, Hopper and Ant (from left to right).

555 **Implementation Details.** For the baseline algorithms, we follow the implementation of BESTDICE
 556 in [57] and the implementations of FQE, MB, IS, DR in [64]. The optimization hyperparameters
 557 including learning rate, optimizer parameter, network architecture, batch size, *etc.*, are selected via
 558 grid search. All the experiments were conducted using V100 GPUs on a multi-node cluster.

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

566 The source code is available at https://anonymous.4open.science/r/spectral_dice-720A.

567 B Primal-Dual Spectral Representation

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

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

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

573 **Corollary 2.** With primal-dual spectral representation (5) where $q(\cdot) \equiv d^{\pi_b}(\cdot)$, under Assumption 1,

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

574 *Proof of Corollary 2.* Recall the primal-dual LP formulation of policy evaluation stated in (2), which
 575 can be equivalently rewritten using the primal-dual spectral representation (5) as follows:

$$\rho(\pi) = \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int d^\pi(s, a) \left[r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] dsda \right\} \quad (10a)$$

$$= \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \int q(s) \pi_b(a|s) \cdot \frac{d^\pi(s, a)}{q(s) \pi_b(a|s)} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')] - Q(s, a) \right] dsda \right\} \quad (10b)$$

$$= \min_{Q(\cdot, \cdot)} \max_{d^\pi(\cdot, \cdot)} \left\{ (1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)] + \mathbb{E}_{s \sim q(\cdot), a \sim \pi_b(\cdot|s), s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[\frac{d^\pi(s, a)}{q(s) \pi_b(\cdot|s)} (r(s, a) + \gamma Q(s', a') - Q(s, a)) \right] \right\} \quad (10c)$$

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

576

577 where in (10b) we perform the IS-style change-of-variable used in DICE estimators (see (3)); in
 578 (10d) we plug in the primal-dual spectral representation of Q^π and d^π stated in (5), as well as the
 579 fact that $q(\cdot) \equiv d^{\pi_b}(\cdot)$. \square

580 B.2 Failure of the Naive Spectral Representation

581 In Section 3.1, it is mentioned that directly applying the naive spectral representation (4) proposed
 582 in Ren et al. [40] induces a complicated representation for $\zeta(\cdot, \cdot)$, which in turn leads to an intractable
 583 optimization (2) for the computation of the DICE estimator. The above point is further elaborated
 584 here in a formal way.

585 Note that, in Lemma 1, the linear representation of Q^π only builds upon the low-rank MDP as-
 586 sumption, and therefore it still holds with the naive spectral representation (4). Meanwhile, it can
 587 be checked that

$$d^\pi(s, a) = \left\langle q(s) \pi(a|s) \mu(s), \underbrace{(1 - \gamma) \omega_0 + \gamma \int d^\pi(\tilde{s}, \tilde{a}) \phi(\tilde{s}, \tilde{a}) d\tilde{s} d\tilde{a}}_{\omega_d^\pi} \right\rangle, \quad (11)$$

588 which can be obtained by plugging the relation $\pi(a|s) \mu(s) = \pi_b(a|s) \mu^\pi(s, a)$ into the linear repre-
 589 sentation of $d^\pi(\cdot, \cdot)$ to eliminate μ^π from the representation. Consequently, the LP formulation (10)
 590 becomes

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

591 which involves an unknown ratio $\frac{\pi(a|s)}{\pi_b(a|s)}$ when the behavior policy π_b is unknown, and is thus in-
 592 tractable.

593 The above failed attempt implies that the policy ratio should be ‘‘absorbed’’ into the representation to
 594 be implicitly learned during representation learning, which exactly inspires the primal-dual spectral
 595 representation (5).

596 B.3 Solving the Minimax Problem via Regularization

597 It is known that directly solving (9) leads to potential numerical instability issues due to the objec-
 598 tive’s linearity in θ_Q and ω_d [19]. Fortunately, it is shown in Yang et al. [57] that certain regulariza-
 599 tion leads to strictly concave inner maximization while keeping the optimal solution ω_d^* unbiased.
 600 Specifically, in practice we may append the following regularizer to the objective in (8):

$$\rho_{\text{reg}}(\pi) = \min_{\theta_Q} \max_{\omega_d} \left\{ (1 - \gamma) \mathbb{E}_{\substack{s \sim \mu_0, \\ a \sim \pi(\cdot|s)}} [\phi(s, a)^\top \theta_Q] + \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), a \sim \pi_b(a|s), \\ s' \sim \mathbb{P}(\cdot|s, a), a' \sim \pi(\cdot|s')}} [(\boldsymbol{\mu}^\pi(s, a)^\top \boldsymbol{\omega}_d) \cdot \right. \\ \left. (r(s, a) + \gamma \phi(s', a')^\top \theta_Q - \phi(s, a)^\top \theta_Q)] - \lambda \mathbb{E}_{(s, a) \sim \mathcal{D}} [f(\hat{\boldsymbol{\mu}}^\pi(s, a)^\top \boldsymbol{\omega}_d)] \right\}. \quad (12)$$

601 Here f is a differentiable convex function with closed and convex Fenchel conjugate f_* (see Ap-
 602 pendix E.1), and $\lambda > 0$ is a tunable constant that controls the magnitude of regularization. It is
 603 evident that the regularized objective is concave in ω_d , which facilitates the inner maximization.
 604 What’s more, it has also been shown that such regularization does not alter the optimal solution ω_d^* ,
 605 as summarized below.

606 **Lemma 5** (Nachum et al. [19], Yang et al. [57]). *The solution $(\theta_Q^{\text{reg},*}, \omega_d^{\text{reg},*})$ to (12) satisfies:*

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

607 where (θ_Q^*, ω_d^*) is the solution to (8).

608 We emphasize that the regularized problem is unbiased only in the sense that $\omega_d^{\text{reg},*} = \omega_d^*$. There-
 609 fore, in general we need to plug $\omega_d^{\text{reg},*}$ back into (8) and solve the outer minimization again to
 610 recover θ_Q^* . Nevertheless, when λ is sufficiently small, we shall regard $\theta_Q^{\text{reg},*} \approx \theta_Q^*$ to relieve the
 611 computational burden.

612 In practice, we can only solve the empirical version of (12), i.e.,

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

613 C Representation Learning Methods and Their Error Bounds

614 In this appendix, we introduce two candidate methods—*ordinary least squares (OLS)* and *noise-*
 615 *contrastive estimation (NCE)*—for the REPLEARN subroutine. Further, we also provide their represen-
 616 tation learning error bounds in the form of Claim 3, which is restated here for readers’ convenience:

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

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

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

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

623 **C.1 Ordinary Least Squares (OLS)**

624 **Method.** Inspired by Ren et al. [40], the objective of OLS can be constructed as follows. Denote
 625 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
 626 behavior policy π_b . Then we plug \mathbb{Q}^π into (5) and rearrange the terms to obtain

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

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

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

628 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,
 629 while the second and third terms can be effectively approximated by sampling from the dataset \mathcal{D}
 630 and the target policy π . Therefore, in practice we learn $(\hat{\phi}, \hat{\boldsymbol{\mu}}^\pi)$ by solving the following optimization:

$$\min_{(\hat{\phi}, \hat{\boldsymbol{\mu}}^\pi) \in \mathcal{F}} \left\{ \widehat{\mathbb{E}}_{(s, a) \sim d^{\pi_b}, (\tilde{s}', \tilde{a}') \sim d^{\pi_b}} \left[(\hat{\phi}(s, a)^\top \hat{\boldsymbol{\mu}}^\pi(\tilde{s}', \tilde{a}'))^2 \right] - 2\widehat{\mathbb{E}}_{(s, a) \sim d^{\pi_b}, (s', a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)} \left[\hat{\phi}(s, a)^\top \hat{\boldsymbol{\mu}}^\pi(s', a') \right] \right\}, \quad (13)$$

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

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

635 **Assumption 4** (regularity for OLS). (1) lower-bounded transition kernel: $\mathbb{P}^\pi(s', a'|s, a) \geq \frac{1}{C_p} > 0$,
 636 $\forall s, a, s', a'$; (2) effective behavior policy coverage: $\frac{d^{\pi_b}(s, a)}{d^{\pi_b}(s', a')} \leq C_{\text{cov}}, \forall s, a, s', a'$.

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

639 **Theorem 6** (OLS learning error). *Under Assumptions 1 to 3 and the additional Assumption 4 for regular-*
 640 *ity, let $(\hat{\phi}, \hat{\boldsymbol{\mu}}^\pi)$ be the solution to (13), and set $\hat{\mathbb{P}}^\pi(s', a'|s, a) := d^{\pi_b}(s', a')\hat{\phi}(s, a)^\top \hat{\boldsymbol{\mu}}^\pi(s')$. Then, for any*
 641 *$\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

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

642 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..

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

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

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

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

646 where we use the fact that $\langle \hat{\phi}(s, a), d^{\pi_b}(s', a')\hat{\boldsymbol{\mu}}^\pi(s', a') \rangle$ is always some valid transition kernel $\tilde{\mathbb{P}}^\pi$
 647 (by Assumption 3), and the additional regularity assumption (Assumption 4). Further, since the

648 family $\tilde{\mathcal{F}}$ is realizable (by Assumption 3), there exists an optimal $f^* \in \tilde{\mathcal{F}}$ such that

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

649 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
650 $1 - \delta$,

$$\int \left(f^*(s, a, s', a') - \hat{f}(s, a, s', a') \right)^2 ds da ds' da' \leq C_{\text{reg}} \cdot \frac{\log(|\mathcal{F}|/\delta)}{N}, \quad (14)$$

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

$$\begin{aligned} & \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[\left\| \mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \\ &= \int d_{\mathbb{P}}^{\pi_b}(s, a) \left| \mathbb{P}^\pi(s', a' | s, a) - \hat{\mathbb{P}}^\pi(s', a' | s, a) \right| ds da ds' da' \end{aligned} \quad (15a)$$

$$= \int \left| \mathbb{Q}^\pi(s', a', s, a) - \hat{\mathbb{Q}}^\pi(s', a', s, a) \right| ds da ds' da' \quad (15b)$$

$$\leq \sqrt{\int \left(\sqrt{\mathbb{Q}^\pi(s', a', s, a)} - \frac{\hat{\mathbb{Q}}^\pi(s', a', s, a)}{\sqrt{\mathbb{Q}^\pi(s', a', s, a)}} \right)^2 ds da ds' da'} \cdot \int \mathbb{Q}^\pi(s', a', s', a') ds da ds' da' \quad (15c)$$

$$= \sqrt{\int \frac{d^{\pi_b}(s', a')}{\mathbb{P}^\pi(s', a' | s, a)} \left(f^*(s, a, s', a') - \hat{f}(s, a, s', a') \right)^2 ds da ds' da'} \quad (15d)$$

$$\leq \sqrt{\max_{s, a, s', a'} \left\{ \frac{d^{\pi_b}(s', a')}{\mathbb{P}^\pi(s', a' | s, a)} \right\}} \cdot \sqrt{C_{\text{reg}} \cdot \frac{\log(|\mathcal{F}|/\delta)}{N}} \quad (15e)$$

$$\leq \sqrt{C_{\mathbb{P}} C_{\text{reg}}} \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}, \quad (15f)$$

653 where in (15b) we use the definition of \mathbb{Q}^π , and define $\hat{\mathbb{Q}}^\pi := d_{\mathbb{P}}^{\pi_b}(s, a) \hat{\mathbb{P}}^\pi(s', a' | s, a)$; in (15c) we use
654 Cauchy-Schwartz inequality; in (15d) we use the definition of \hat{f} and f^* ; in (15e) we plug in the PAC
655 bound (14); in (15f) we use Assumption 4 to bound the coefficient. This completes the proof. \square

656 C.2 Noise-Contrastive Learning (NCE)

657 **Method.** NCE is a widely used method for contrastive representation learning in RL [47, 48]. To
658 learn $(\hat{\phi}, \hat{\mu}^\pi)$, we consider a binary contrastive learning objective [47]:

$$\min_{(\hat{\phi}, \hat{\mu}^\pi) \in \mathcal{F}} \left[\mathbb{E}_{(s,a) \sim d^{\pi_b}} \left[\mathbb{E}_{(s',a') \sim \mathbb{P}^\pi(\cdot, \cdot | s, a)} \left[\log \left(1 + \frac{1}{\hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a')} \right) \right] \right] + \mathbb{E}_{(s',a') \sim P_{\text{neg}}} \left[\log \left(1 + \hat{\phi}(s, a)^\top \hat{\mu}^\pi(s', a') \right) \right] \right], \quad (16)$$

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

662 The following derivations follow a similar pathway as those in [47]. For notational consistency that
663 facilitates the application of known results, we introduce the following auxiliary notations. Define

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

664 For clarity, we augment the sampled transitions to include a label y indicating whether the sample
665 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
666 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
667 define the augmented dataset

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

668 In this way, the NCE objective (16) can be equivalently rewritten (in MLE format) as

$$\max_{f \in \tilde{\mathcal{F}}} \widehat{\mathbb{E}}_{(s,a,s',a',y) \sim d^{\tilde{\mathcal{D}}}} [\log \psi_f(s, a, s', a', y)], \quad (17)$$

669 where the likelihood function ψ_f is defined by

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

670 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
 671 f^* that optimizes the unconstrained non-empirical version of (17), ψ_{f^*} can be interpreted as the
 672 probability of obtaining label y given (s, a, s', a') , as summarized in the following lemma that is
 673 similar to Lemma C.1 in Qiu et al. [47].

674 **Lemma 7** (non-empirical solution to NCE). *The optimal solution to the unconstrained non-empirical*
 675 *version of (17), i.e., $f^* := \max_f \mathbb{E}_{(s,a,s',a',y) \sim d^{\tilde{\mathcal{D}}}} [\log \psi_f(s, a, s', a', y)]$, is characterized by*

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

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

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

677 Here $\mathbb{H}(\cdot; \cdot)$ is the cross entropy between distributions, which, by Gibbs' inequality, is minimized
 678 only when

$$\Pr(y | s, a, s', a') = \psi_{f^*}(s, a, s', a', y) = \left(\frac{f^*(s, a, s', a')}{1 + f^*(s, a, s', a')} \right)^y \cdot \left(\frac{1}{1 + f^*(s, a, s', a')} \right)^{1-y}. \quad (18)$$

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

$$\Pr(y = 1 | s, a, s', a') = \frac{\Pr(s', a' | s, a, y = 1) \Pr(y = 1 | s, a)}{\sum_{y \in \mathcal{Y}} \Pr(s', a' | s, a, y) \Pr(y | s, a)} = \frac{\mathbb{P}^\pi(s', a' | s, a)}{P_{\text{neg}}(s', a') + \mathbb{P}^\pi(s', a' | s, a)}. \quad (19)$$

680 Comparing (18) and (19) gives

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

681 This completes the proof. \square

682 *Remark 3.* For conciseness, here we slightly abuse the notation $\Pr(\cdot)$ to denote the distribution (den-
 683 sity or mass) of joint and conditional distributions involving random variables $(s, a, s', a', y) \sim d^{\tilde{\mathcal{D}}}$.
 684 Specifically, we write $\Pr(\dots, x, \dots)$ to indicate an arbitrary value x taken by the random variable,
 685 and we also write $\Pr(\dots, x = x_0, \dots)$ to emphasize the value x_0 taken by that random variable.

686 Lemma 7 is important in that it echoes the form of primal-dual spectral representation in (5). Specif-
 687 ically, we shall take $P_{\text{neg}}(\cdot, \cdot) \equiv d^{\pi_b}(\cdot, \cdot)$ for an exact match, which is also implementable using offline
 688 data since d^{π_b} can be effectively approximated by sampling the trajectories. We will stick to this
 689 choice of P_{neg} from now on.

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

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

695 We point out that Assumption 5 is a standard assumption for the negative sampling distribu-
 696 tion [47], aiming at eliminating the cases where certain transitions are scarcely drawn as nega-
 697 tive samples and thus obstruct efficient representation learning for those cases. The assumption
 698 is also slightly stronger than the effective behavior policy coverage assumption required by the OLS
 699 method (see Assumption 4).

700 **Theorem 8** (NCE learning error). *Under Assumptions 1 to 3 and the additional Assumption 5 for*
 701 *regularity, let $(\hat{\phi}, \hat{\mu}^\pi)$ be the solution to (16) with $P_{\text{neg}}(\cdot, \cdot) \equiv d^{\pi_b}(\cdot, \cdot)$, and set $\hat{\mathbb{P}}^\pi(s', a'|s, a) :=$
 702 $d^{\pi_b}(s', a')\hat{\phi}(s, a)^\top \hat{\mu}^\pi(s')$. Then, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

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

703 The proof of Theorem 8 largely follows the same pathway and techniques established in Qiu et al.
 704 [47]. Nevertheless, our proof is less technically involved since the offline non-episodic setting sig-
 705 nificantly weakens the correlation between samples. For the sake of completeness, we restate the
 706 proof below.

707 *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
 708 in turn be calculated using Bayes' rule as follows:

$$\begin{aligned} \Pr(s', a' | s, a) &= \Pr(s', a' | s, a, y = 0)\Pr(y = 0 | s, a) + \Pr(s', a' | s, a, y = 1)\Pr(y = 1 | s, a) \\ &= \frac{1}{2}(\mathbb{P}^\pi(s', a' | s, a) + P_{\text{neg}}(s', a')). \end{aligned} \quad (20)$$

709 Here we use the fact that the data distribution $d^{\tilde{\mathcal{D}}}$ implicitly assigns an equal number of labels as
 710 $y = 0$ and $y = 1$ by the design of NCE objective (16). Since $\Pr(s', a' | s, a)$ is a constant that is
 711 independent from f , we can further rewrite the NCE objective to be

$$\arg \max_{f \in \tilde{\mathcal{F}}} \left\{ \hat{\mathbb{E}}_{(s,a,s',a',y) \in \tilde{\mathcal{D}}} [\log \Pr_f(y | s, a, s', a')] \right\} = \arg \max_{f \in \tilde{\mathcal{F}}} \left\{ \hat{\mathbb{E}}_{(s,a,s',a',y) \in \tilde{\mathcal{D}}} [\log \Pr_f(y, s', a' | s, a)] \right\}, \quad (21)$$

712 where we define the shorthand notations

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

713 for any $f \in \tilde{\mathcal{F}}$. Note that the right-hand side of (21) is in the desired MLE form, with ground-
 714 truth conditional density $\Pr_{f^*}(y, s', a' | s, a)$ specified by some $f^* \in \tilde{\mathcal{F}}$, thanks to the realizability
 715 assumption (Assumption 3). Now, using the PAC bound for MLE shown in Agarwal et al. [41] (see
 716 Lemma 19), we have

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

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

$$\mathbb{E}_{(s,a) \sim d_{\mathbb{P}^b}^{\pi_b}} \left[\left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1^2 \right] \leq \frac{8 \log(|\mathcal{F}|/\delta)}{N}. \quad (22)$$

719 We proceed to further relate (22) with the desired format. For this purpose, note that

$$\begin{aligned} &\left\| \Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a) \right\|_1 \\ &= \left\| \Pr_{\hat{f}}(y = 1, \cdot, \cdot | s, a) - \Pr_{f^*}(y = 1, \cdot, \cdot | s, a) \right\|_1 + \left\| \Pr_{\hat{f}}(y = 0, \cdot, \cdot | s, a) - \Pr_{f^*}(y = 0, \cdot, \cdot | s, a) \right\|_1 \end{aligned} \quad (23a)$$

$$= 2 \left\| \frac{\Pr(\cdot, \cdot | s, a)}{1 + \hat{f}(s, a, \cdot, \cdot)} - \frac{\Pr(\cdot, \cdot | s, a)}{1 + f^*(s, a, \cdot, \cdot)} \right\|_1 \quad (23b)$$

$$= 2 \int \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} ds' da', \quad (23c)$$

720 where in (23a) we use the definition of L^1 -norm; in (23b) we use the fact that

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

721 Now, plugging Lemma 7 and (20) into the integrand in (23c), we have

$$\begin{aligned} & \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} \\ &= \frac{|\mathbb{P}^\pi(s', a' | s, a) / P_{\text{neg}}(s', a') - \hat{f}(s, a, s', a')| \cdot \frac{1}{2} (\mathbb{P}^\pi(s', a' | s, a) + P_{\text{neg}}(s', a'))}{(1 + \hat{f}(s, a, s', a'))(1 + \mathbb{P}^\pi(s', a' | s, a) / P_{\text{neg}}(s', a'))} \end{aligned} \quad (24a)$$

$$= \frac{|\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a') \hat{f}(s, a, s', a')|}{2(1 + \hat{f}(s, a, s', a'))} \quad (24b)$$

$$\geq \frac{|\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a') \hat{f}(s, a, s', a')|}{2(1 + C_d)}, \quad (24c)$$

722 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 (24c). Hence,

$$\begin{aligned} & \|\mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a)\|_1 \\ &= \int |\mathbb{P}^\pi(s', a' | s, a) - P_{\text{neg}}(s', a') \hat{f}(s, a, s', a')| ds' da' \end{aligned} \quad (25a)$$

$$\leq 2(1 + C_d) \int \frac{|\hat{f}(s, a, s', a') - f^*(s, a, s', a')| \cdot \Pr(s', a' | s, a)}{(1 + \hat{f}(s, a, s', a'))(1 + f^*(s, a, s', a'))} ds' da' \quad (25b)$$

$$= (1 + C_d) \|\Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a)\|_1, \quad (25c)$$

723 where we use $\hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) = P_{\text{neg}}(\cdot, \cdot) \hat{f}(s, a, \cdot, \cdot)$ in (25a), (24) in (25b), and (23) in (25c). Finally,

$$\begin{aligned} & \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[\|\mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a)\|_1 \right] \\ & \leq \sqrt{\mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[\|\mathbb{P}^\pi(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a)\|_1^2 \right]} \end{aligned} \quad (26a)$$

$$\leq \sqrt{(1 + C_d)^2 \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi_b}} \left[\|\Pr_{\hat{f}}(\cdot, \cdot, \cdot | s, a) - \Pr_{f^*}(\cdot, \cdot, \cdot | s, a)\|_1^2 \right]} \quad (26b)$$

$$\leq 2\sqrt{2}(1 + C_d) \cdot \sqrt{\frac{\log(|\mathcal{F}|/\delta)}{N}}, \quad (26c)$$

724 where we use Cauchy-Schwartz inequality in (26a), (25) in (26b), and (22) in (26c). \square

725 D Sample Complexity Guarantee

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

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

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

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

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

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

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

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

$$\begin{aligned} & \rho_{\hat{\mathbb{P}}}(\pi) - \rho_{\mathbb{P}}(\pi) \\ &= \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_{\hat{\mathbb{P}}}^\pi} \left[\mathbb{E}_{s' \sim \hat{\mathbb{P}}(\cdot|s,a)} [V_{\hat{\mathbb{P}}}^\pi(s')] - \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)} [V_{\hat{\mathbb{P}}}^\pi(s')] \right] \end{aligned} \quad (27a)$$

$$= \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_{\hat{\mathbb{P}}}^\pi} \left[\mathbb{E}_{(s',a') \sim \hat{\mathbb{P}}^\pi(\cdot, \cdot|s,a)} [Q_{\hat{\mathbb{P}}}^\pi(s', a')] - \mathbb{E}_{(s',a') \sim \mathbb{P}^\pi(\cdot, \cdot|s,a)} [Q_{\hat{\mathbb{P}}}^\pi(s', a')] \right] \quad (27b)$$

$$= \frac{\gamma}{1-\gamma} \int d_{\hat{\mathbb{P}}}^\pi(s, a) ds da \int Q_{\hat{\mathbb{P}}}^\pi(s', a') \left(\hat{\mathbb{P}}^\pi(s', a'|s, a) - \mathbb{P}^\pi(s', a'|s, a) \right) ds' da' \quad (27c)$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \int C_\infty^\pi d_{\hat{\mathbb{P}}}^{\pi_b}(s, a) ds da \int |\mathbb{P}^\pi(s', a'|s, a) - \hat{\mathbb{P}}^\pi(s', a'|s, a)| ds' da' \quad (27d)$$

$$= \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \mathbb{E}_{(s,a) \sim d_{\hat{\mathbb{P}}}^{\pi_b}} \left[\|\mathbb{P}^\pi(s', a'|s, a) - \hat{\mathbb{P}}^\pi(s', a'|s, a)\|_1 \right] \quad (27e)$$

$$\leq \frac{\gamma C_\infty^\pi}{(1-\gamma)^2} \cdot \xi(|\mathcal{F}|, N, \delta), \quad (27f)$$

742 where in (27a) we use the Simulation Lemma; in (27b) we use the relationship between value func-
 743 tions; in (27d) we plug in $d_{\hat{\mathbb{P}}}^\pi(s, a) \leq C_\infty^\pi d_{\hat{\mathbb{P}}}^{\pi_b}(s, a)$ (Assumption 2) and the fact that $Q_{\hat{\mathbb{P}}}^\pi(\cdot, \cdot) \leq \frac{1}{1-\gamma}$;
 744 in (27f) we use Claim 3. \square

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

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

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

$$\min_{\mathbf{x}} \max_{\mathbf{y}} F_1(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x}} \left\{ \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \max_{\mathbf{y}} \{F_1(\mathbf{x}, \mathbf{y}) - F_2(\mathbf{x}, \mathbf{y})\} \right\} \quad (28a)$$

$$\leq \min_{\mathbf{x}} \left\{ \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \varepsilon \right\} \quad (28b)$$

$$= \min_{\mathbf{x}} \max_{\mathbf{y}} F_2(\mathbf{x}, \mathbf{y}) + \varepsilon, \quad (28c)$$

750 where in (28a) 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

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

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

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

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

$$\bar{\rho}(\pi) = \min_{\boldsymbol{\theta}_Q} \max_{\boldsymbol{\omega}_d} \bar{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d), \quad \rho_{\hat{\mathbb{P}}}(\pi) = \min_{\boldsymbol{\theta}_Q} \max_{\boldsymbol{\omega}_d} \hat{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d).$$

754 Then we can show that

$$|\bar{\rho}(\pi) - \rho_{\hat{\mathbb{P}}}(\pi)| = \left| \min_{\boldsymbol{\theta}_Q} \max_{\boldsymbol{\omega}_d} \bar{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d) - \min_{\boldsymbol{\theta}_Q} \max_{\boldsymbol{\omega}_d} \hat{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d) \right| \leq \left| \bar{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d) - \hat{F}(\boldsymbol{\theta}_Q, \boldsymbol{\omega}_d) \right| \quad (29a)$$

$$= \left| \int d_{\hat{\mathbb{P}}^{\pi_b}}(s, a) \left(\mathbb{P}^{\pi}(s', a' | s, a) - \hat{\mathbb{P}}^{\pi}(s', a' | s, a) \right) (\hat{\boldsymbol{\mu}}^{\pi}(s, a)^{\top} \boldsymbol{\omega}_d) \cdot \right. \\ \left. \left(r(s, a) + \gamma \hat{\boldsymbol{\phi}}(s', a')^{\top} \boldsymbol{\theta}_Q - \hat{\boldsymbol{\phi}}(s, a)^{\top} \boldsymbol{\theta}_Q \right) ds da ds' da' \right| \quad (29b)$$

$$\leq \int d_{\hat{\mathbb{P}}^{\pi_b}}(s, a) \left| \mathbb{P}^{\pi}(s', a' | s, a) - \hat{\mathbb{P}}^{\pi}(s', a' | s, a) \right| \cdot |\hat{\boldsymbol{\mu}}^{\pi}(s, a)^{\top} \boldsymbol{\omega}_d| \cdot \\ \left| r(s, a) + \gamma \hat{\boldsymbol{\phi}}(s', a')^{\top} \boldsymbol{\theta}_Q - \hat{\boldsymbol{\phi}}(s, a)^{\top} \boldsymbol{\theta}_Q \right| ds da ds' da' \quad (29c)$$

$$\leq \mathbb{E}_{(s, a) \sim d_{\hat{\mathbb{P}}^{\pi_b}}} \left[\left\| \mathbb{P}^{\pi}(\cdot, \cdot | s, a) - \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a) \right\|_1 \cdot C_{\infty}^{\pi} \cdot \frac{1}{1-\gamma} \right] \quad (29d)$$

$$= \frac{C_{\infty}^{\pi}}{1-\gamma} \cdot \xi(|\mathcal{F}|, N, \delta), \quad (29e)$$

755 where in (29a) we use Lemma 10; in (29c) we use the integral triangle inequality; in (29d) we plug
756 in $|\hat{\boldsymbol{\mu}}^{\pi}(s, a)^{\top} \boldsymbol{\omega}_d| \leq C_{\infty}^{\pi}$ and $|r(s, a) + \gamma \hat{\boldsymbol{\phi}}(s', a')^{\top} \boldsymbol{\theta}_Q - \hat{\boldsymbol{\phi}}(s, a)^{\top} \boldsymbol{\theta}_Q| \leq \frac{1}{1-\gamma}$ (see Remark 2); in (29e)
757 we use Claim 3. \square

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

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

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

761 *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

$$F(s, a, s', a') := (\hat{\boldsymbol{\mu}}^{\pi}(s, a)^{\top} \boldsymbol{\omega}_d) (r(s, a) + \gamma \hat{\boldsymbol{\phi}}(s', a')^{\top} \boldsymbol{\theta}_Q - \hat{\boldsymbol{\phi}}(s, a)^{\top} \boldsymbol{\theta}_Q).$$

762 Note that $|\hat{\boldsymbol{\mu}}^{\pi}(s, a)^{\top} \boldsymbol{\omega}_d| \leq C_{\infty}^{\pi}$ and $|r(s, a) + \gamma \hat{\boldsymbol{\phi}}(s', a')^{\top} \boldsymbol{\theta}_Q - \hat{\boldsymbol{\phi}}(s, a)^{\top} \boldsymbol{\theta}_Q| \leq \frac{1}{1-\gamma}$ (see Remark 2), we
763 have

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

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

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

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

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

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

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

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

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

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

$$\begin{aligned} \mathcal{C}_1 &: \mathbb{E}_{(s,a) \sim d^{\pi_b}} \left[\left\| \hat{\mathbb{P}}^\pi(\cdot, \cdot | s, a) - \mathbb{P}^\pi(\cdot, \cdot | s, a) \right\|_1 \right] \leq \xi(|\mathcal{F}|, N, \delta/2), \\ \mathcal{C}_2 &: \left| \hat{\mathbb{E}}_{\substack{(s,a,s') \sim \mathcal{D}, \\ a' \sim \pi(\cdot | s')}} [F(s, a, s', a')] - \mathbb{E}_{\substack{s \sim d^{\pi_b}(\cdot), \\ s' \sim \mathbb{P}(\cdot | s, a), \\ a' \sim \pi(\cdot | s')}} [F(s, a, s', a')] \right| \leq \frac{C_\infty^\pi}{1-\gamma} \sqrt{\frac{\log(1/\delta)}{2N}}. \end{aligned}$$

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

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

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

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

774 This completes the proof. \square

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

778 **Corollary 13** (sample complexity of OLS-based SPECTRALDICE). *Under Assumptions 1 to 3 and the
779 additional Assumption 4 for regularity, let $(\hat{\phi}, \hat{\mu}^\pi)$ be the solution to the OLS problem (13). Then, for any
780 $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

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

781 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..

782 **Corollary 14** (sample complexity of NCE-based SPECTRALDICE). *Under Assumptions 1 to 3 and the
783 additional Assumption 5 for regularity, let $(\hat{\phi}, \hat{\mu}^\pi)$ be the solution to the NCE problem (16) with $P_{\text{neg}}(\cdot, \cdot) \equiv$
784 $d^{\pi_b}(\cdot, \cdot)$. Then, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have*

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

785 *Remark 4* (Sampling the dataset). Throughout this paper, we have been slightly abusing the notation
786 $(s, a, s') \sim \mathcal{D}$, which is a little subtle in practice since only trajectories (rather than transitions) are
787 collected. To ensure the correct data distribution $d^{\mathcal{D}}(s, a) = d^{\pi_b}(s, a)$, we shall first randomly sample
788 the trajectories, within which we sample each transition $(s_t, a_t, s_{t+1}, a_{t+1})$ with probability $(1-\gamma)\gamma^t$.

789 E Technical Lemmas

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

791 E.1 f -Divergence

792 **Definition 1** (f -divergence). Let \mathbb{P} and \mathbb{Q} be two probabilities distribution over a sample space \mathcal{X} ,
793 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
794 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

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

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

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

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

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

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

799 *Proof.* See Theorem 4.4 in Broniatowski and Keziou [65]. □

800 E.2 Concentration Inequalities

801 **Lemma 16** (Hoeffding's inequality, Hoeffding [66]). *Let X_1, X_2, \dots, X_N be i.i.d. random variables*
 802 *with mean μ and taking values in $[a, b]$ almost surely. Then for any $\varepsilon > 0$ we have*

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

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

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

804 **Lemma 17** (Bernstein's inequality, Bernstein [67]). *Let X_1, X_2, \dots, X_N be i.i.d. random variables with*
 805 *mean μ , variance σ^2 , and bounded range $|X_i - \mu| \leq B$ almost surely. Then with probability at least $1 - \delta$,*
 806 *we have*

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

807 E.3 Statistical Learning: PAC Bounds

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

811 **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*

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

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

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

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

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

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

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

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

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

$$\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[Z_i(f)] = \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - Y_i)^2 - (f^*(X_i) - Y_i)^2] \quad (30a)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}\left[\left((f(X_i) - f^*(X_i)) + (f^*(X_i) - Y_i)\right)^2 - (f^*(X_i) - Y_i)^2\right] \quad (30b)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - f^*(X_i))^2] + 2\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - f^*(X_i))(f^*(X_i) - Y_i)] \quad (30c)$$

$$= \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - f^*(X_i))^2] =: \mathcal{E}(f), \quad (30d)$$

820 where in (30c) we use the following fact $\mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - f^*(X_i))(f^*(X_i) - Y_i)] = \mathbb{E}_{X_i}[(f(X_i) -$
821 $f^*(X_i)) \cdot \mathbb{E}_{Y_i \sim \mathcal{M}(\cdot|X_i)}[f^*(X_i) - Y_i]] = 0$ that directly follows from the definition of f^* . Similarly, for
822 any $t \in [T]$,

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

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

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

$$\leq 4(b-a)^2 \mathbb{E}_{(X_i, Y_i) \sim \mathcal{M}}[(f(X_i) - f^*(X_i))^2] = 4(b-a)^2 \mathcal{E}(f). \quad (31d)$$

823 where in (31a) we simply drop the second term, and in (31c) we use the fact $f(X_i) + f^*(X_i) - 2Y_i \in$
824 $[-2(b-a), 2(b-a)]$ 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
825 $f(x) - y \in [-(b-a), b-a]$, implying $Z_i(f) \in [-(b-a), b-a]$ and $\mathbb{E}[Z_i(f)] \in [-(b-a), (b-a)]$.
826 Therefore, $|Z_i(f) - \mathbb{E}[Z_i(f)]| \leq 2(b-a)$. Then by Bernstein's inequality (Lemma 17), we conclude
827 that, with probability at least $1 - \delta$,

$$\mathbb{E}[Z_i(f)] - \frac{1}{N} \sum_{i=1}^N Z_i(f) \leq \sqrt{\text{Var}[Z_i(f)]} \sqrt{\frac{2 \log(1/\delta)}{N}} + \frac{2(b-a) \log(1/\delta)}{3N}. \quad (32)$$

828 To proceed, plug (30) and (31) into (32), and we have

$$\mathcal{E}(f) - (\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) \leq 2(b-a) \sqrt{\mathcal{E}(f)} \sqrt{\frac{2 \log(1/\delta)}{N}} + \frac{2(b-a) \log(1/\delta)}{3N} \quad (33a)$$

$$\leq \left(\frac{1}{2} \mathcal{E}(f) + \frac{4(b-a)^2 \log(1/\delta)}{N}\right) + \frac{2(b-a) \log(1/\delta)}{3N}, \quad (33b)$$

829 where in (33a) we apply the AM-GM inequality. Finally, we rearrange the terms to obtain

$$\mathcal{E}(f) \leq 2(\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) + \frac{C_{\text{reg}} \log(1/\delta)}{N} \quad (34)$$

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

$$\mathcal{E}(f) \leq 2(\hat{\mathcal{L}}(f) - \hat{\mathcal{L}}(f^*)) + \frac{C_{\text{reg}} \log(|\mathcal{F}|/\delta)}{N}, \quad \forall f \in \mathcal{F}. \quad (35)$$

832 In particular, (35) also applies to the ERM estimator \hat{f} , which gives

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

833 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.
834 \square

835 **Lemma 19** (PAC bound for MLE, Agarwal et al. [41]). Consider a conditional probability estimation
836 problem over a finite family $\mathcal{F} = \{f : (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}\}$, where the objective is to estimate $f^*(x, y) := \mathbb{P}(y|x)$.
837 Suppose the ground truth $f^* \in \mathcal{F}$ (realizability), and we have access to (potentially correlated) samples
838 $\{(x_i, y_i) \mid i \in [N]\}$ 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
839 process) and $y_i \sim \mathbb{P}(\cdot|x_i)$. Let the Maximum Likelihood Estimator (MLE) be

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

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

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

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

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

842 E.4 Simulation Lemma in MDPs

843 The following Simulation Lemma is a simplified version of Lemma 21 in Uehara et al. [42].

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

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

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

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

846 Therefore, since $\rho_{\mathbb{P}}(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\mathbb{P}}^{\pi}} [r(s, a)]$ 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

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

847 where in the last equality we plug in the Bellman equation

$$Q_{\hat{\mathbb{P}}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{(s',a') \sim \hat{\mathbb{P}}^{\pi}(\cdot, \cdot | s, a)} [Q_{\hat{\mathbb{P}}}^{\pi}(s', a')].$$

848 Finally, we leverage the relationship between Q - and V -functions to complete the proof. \square