
DualDICE: Behavior-Agnostic Estimation of Discounted Stationary Distribution Corrections

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In many real-world reinforcement learning applications, access to the environ-
2 ment is limited to a fixed dataset, instead of direct (online) interaction with the
3 environment. When using this data for either evaluation or training of a new pol-
4 icy, accurate estimates of *discounted stationary distribution ratios* — correction
5 terms which quantify the likelihood that the new policy will experience a certain
6 state-action pair normalized by the probability with which the state-action pair
7 appears in the dataset — can improve accuracy and performance. In this work,
8 we propose an algorithm, DualDICE, for estimating these quantities. In contrast
9 to previous approaches, our algorithm is agnostic to knowledge of the behavior
10 policy (or policies) used to generate the dataset. Furthermore, it eschews any
11 direct use of importance weights, thus avoiding potential optimization instabilities
12 endemic of previous methods. In addition to providing theoretical guarantees, we
13 present an empirical study of our algorithm applied to off-policy policy evaluation
14 and find that our algorithm significantly improves accuracy compared to existing
15 techniques.

16 1 Introduction

17 Reinforcement learning (RL) has recently demonstrated a number of successes in various domains,
18 such as games [31], robotics [1], and conversational systems [15, 24]. These successes have often
19 hinged on the use of simulators to provide large amounts of experience necessary for RL algorithms.
20 While this is reasonable in game environments, where the game is often a simulator itself, and some
21 simple real-world tasks can be simulated to an accurate enough degree, in general one does not have
22 such direct or easy access to the environment. Furthermore, in many real-world domains such as
23 medicine [32], recommendation [25], and education [30], the deployment of a new policy, even just
24 for the sake of performance evaluation, may be expensive and risky. In these applications, access
25 to the environment is usually in the form of *off-policy* data [46], logged experience collected by
26 potentially multiple and possibly unknown *behavior* policies.

27 State-of-the-art methods which consider this more realistic setting — either for policy evaluation
28 or policy improvement — often rely on estimating (*discounted*) *stationary distribution ratios* or
29 *corrections*. For each state and action in the environment, these quantities measure the likelihood
30 that one’s current *target* policy will experience the state-action pair normalized by the probability
31 with which the state-action pair appears in the off-policy data. Proper estimation of these ratios can
32 improve the accuracy of policy evaluation [27] and the stability of policy learning [16, 18, 28, 47]. In
33 general, these ratios are difficult to compute, let alone estimate, as they rely not only on the probability
34 that the target policy will take the desired action at the relevant state, but also on the probability that
35 the target policy’s interactions with the environment dynamics will lead it to the relevant state.

36 Several methods to estimate these ratios have been proposed recently [16, 18, 27], all based on the
37 steady-state property of stationary distributions of Markov processes [19]. This property may be

expressed locally with respect to state-action-next-state tuples, and is therefore amenable to stochastic optimization algorithms. However, these methods possess several issues when applied in practice: First, these methods require knowledge of the probability distribution used for each sampled action appearing in the off-policy data. In practice, these probabilities are usually not known and difficult to estimate, especially in the case of multiple, non-Markovian behavior policies. Second, the loss functions of these algorithms involve per-step importance ratios (the ratio of action sample probability with respect to the target policy versus the behavior policy). Depending on how far the behavior policy is from the target policy, these quantities may have large variance, and thus have a detrimental effect on stochastic optimization algorithms.

In this work, we propose *Dual stationary DIstribution Correction Estimation (DualDICE)*, a new method for estimating discounted stationary distribution ratios. It is agnostic to the number or type of behavior policies used for collecting the off-policy data. Moreover, the objective function of our algorithm does not involve any per-step importance ratios, and so our solution is less likely to be affected by their high variance. We provide theoretical guarantees on the convergence of our algorithm and evaluate it on a number of off-policy policy evaluation benchmarks. We find that DualDICE can consistently, and often significantly, improve performance compared to previous algorithms for estimating stationary distribution ratios.

2 Background

We consider a Markov Decision Process (MDP) setting [39], in which the environment is specified by a tuple $\mathcal{M} = \langle S, A, R, T, \beta \rangle$, consisting of a state space, an action space, a reward function, a transition probability function, and an initial state distribution. A policy π interacts with the environment iteratively, starting with an initial state $s_0 \sim \beta$. At step $t = 0, 1, \dots$, the policy produces a distribution $\pi(\cdot|s_t)$ over the actions A , from which an action a_t is sampled and applied to the environment. The environment stochastically produces a scalar reward $r_t \sim R(s_t, a_t)$ and a next state $s_{t+1} \sim T(s_t, a_t)$. In this work, we consider infinite-horizon environments and the γ -discounted reward criterion for $\gamma \in [0, 1)$. It is clear that any finite-horizon environment may be interpreted as infinite-horizon by considering an augmented state space with an extra terminal state which continually loops onto itself with zero reward.

2.1 Off-Policy Policy Evaluation

Given a *target* policy π , we are interested in estimating its value, defined as the normalized expected per-step reward obtained by following the policy:

$$\rho(\pi) := (1 - \gamma) \cdot \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 \sim \beta, \forall t, a_t \sim \pi(s_t), r_t \sim R(s_t, a_t), s_{t+1} \sim T(s_t, a_t) \right]. \quad (1)$$

The off-policy policy evaluation (OPE) problem studied here is to estimate $\rho(\pi)$ using a fixed set \mathcal{D} of transitions (s, a, r, s') sampled in a certain way. This is a very general scenario: \mathcal{D} can be collected by a single behavior policy (as in most previous work), multiple behavior policies, or an oracle sampler, among others. In the special case where \mathcal{D} contains entire trajectories collected by a known behavior policy μ , one may use *importance sampling* (IS) to estimate $\rho(\pi)$. Specifically, given a finite-length trajectory $\tau = (s_0, a_0, r_0, \dots, s_H)$ collected by μ , the IS estimate of ρ based on τ is estimated by [38]: $(1 - \gamma) \left(\prod_{t=0}^{H-1} \frac{\pi(a_t|s_t)}{\mu(a_t|s_t)} \right) \left(\sum_{t=0}^{H-1} \gamma^t r_t \right)$. Although many improvements exist [e.g., 13, 21, 38, 50], importance-weighting the entire trajectory can suffer from exponentially high variance, which is known as “the curse of horizon” [26, 27].

To avoid exponential dependence on trajectory length, one may weight the states by their *long-term* occupancy measure. First, observe that the policy value may be re-expressed as,

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi, r \sim R(s,a)} [r],$$

where

$$d^\pi(s, a) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s, a_t = a \mid s_0 \sim \beta, \forall t, a_t \sim \pi(s_t), s_{t+1} \sim T(s_t, a_t)), \quad (2)$$

is the *normalized discounted stationary distribution* over state-actions with respect to π . One may define the discounted stationary distribution over states analogously, and we slightly abuse notation by denoting it as $d^\pi(s)$; note that $d^\pi(s, a) = d^\pi(s) \pi(a|s)$. If \mathcal{D} consists of trajectories collected by a behavior policy μ , then the policy value may be estimated as,

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\mu, r \sim R(s,a)} [w_{\pi/\mu}(s, a) \cdot r],$$

85 where $w_{\pi/\mu}(s, a) = \frac{d^\pi(s, a)}{d^\mu(s, a)}$ is the *discounted stationary distribution correction*. The key challenge
 86 is in estimating these correction terms using data drawn from d^μ .

87 2.2 Learning Stationary Distribution Corrections

88 We provide a brief summary of previous methods for estimating the stationary distribution corrections.
 89 The ones that are most relevant to our work are a suite of recent techniques [16, 18, 27], which are all
 90 essentially based on the following steady-state property of stationary Markov processes:

$$d^\pi(s') = (1 - \gamma)\beta(s') + \gamma \sum_{s \in S} \sum_{a \in A} d^\pi(s) \pi(a|s) T(s'|s, a), \quad \forall s' \in S, \quad (3)$$

91 where we have simplified the identity by restricting to discrete state and action spaces. This identity
 92 simply reflects the conservation of flow of the stationary distribution: At each timestep, the flow out
 93 of s' (the LHS) must equal the flow into s' (the RHS). Given a behavior policy μ , equation 3 can be
 94 equivalently rewritten in terms of the stationary distribution corrections, i.e., for any given $s' \in S$,

$$\mathbb{E}_{(s_t, a_t, s_{t+1}) \sim d^\mu} [\text{TD}(s_t, a_t, s_{t+1} | w_{\pi/\mu}) | s_{t+1} = s'] = 0, \quad (4)$$

95 where

$$\text{TD}(s, a, s' | w_{\pi/\mu}) := -w_{\pi/\mu}(s') + (1 - \gamma)\beta(s') + \gamma w_{\pi/\mu}(s) \cdot \frac{\pi(a|s)}{\mu(a|s)},$$

96 provided that $\mu(a|s) > 0$ whenever $\pi(a|s) > 0$. The quantity TD can be viewed as a *temporal differ-*
 97 *ence* associated with $w_{\pi/\mu}$. Accordingly, previous works optimize loss functions which minimize
 98 this TD error using samples from d^μ . We emphasize that although $w_{\pi/\mu}$ is associated with a temporal
 99 difference, it does not satisfy a Bellman recurrence in the usual sense [3]. Indeed, note that equation 3
 100 is written “backwards”: The occupancy measure of a state s' is written as a (discounted) function of
 101 *previous* states, as opposed to vice-versa. This will serve as a key differentiator between our algorithm
 102 and these previous methods.

103 2.3 Off-Policy Estimation with Multiple Unknown Behavior Policies

104 While the previous algorithms are promising, they have several limitations when applied in practice:

- 105 • The off-policy experience distribution d^μ is with respect to a single, Markovian behavior policy μ ,
 106 and this policy must be known during optimization. In practice, off-policy data often comes from
 107 multiple, unknown behavior policies.
- 108 • Computing the TD error in equation 4 requires the use of per-step importance ratios
 109 $\pi(a_t|s_t)/\mu(a_t|s_t)$ at every state-action sample (s_t, a_t) . Depending on how far the behavior policy
 110 is from the target policy, these quantities may have high variance, which can have a detrimental
 111 effect on the convergence of any stochastic optimization algorithm that is used to estimate $w_{\pi/\mu}$.

112 The method we derive below will be free of the aforementioned issues, avoiding unnecessary
 113 requirements on the form of the off-policy data collection as well as explicit uses of importance
 114 ratios. Rather, we consider the general setting where \mathcal{D} consists of *transitions* sampled in an unknown
 115 fashion. Since \mathcal{D} contains rewards and next states, we will often slightly abuse notation and write not
 116 only $(s, a) \sim d^\mathcal{D}$ but also $(s, a, r) \sim d^\mathcal{D}$ and $(s, a, s') \sim d^\mathcal{D}$, where the notation $d^\mathcal{D}$ emphasizes that,
 117 unlike previously, \mathcal{D} is not the result of a single, known behavior policy. The target policy’s value
 118 can be equivalently written as,

$$\rho(\pi) = \mathbb{E}_{(s, a, r) \sim d^\mathcal{D}} [w_{\pi/\mathcal{D}}(s, a) \cdot r], \quad (5)$$

119 where the correction terms are given by $w_{\pi/\mathcal{D}}(s, a) := \frac{d^\pi(s, a)}{d^\mathcal{D}(s, a)}$, and our algorithm will focus on
 120 estimating these correction terms. Rather than relying on the assumption that \mathcal{D} is the result of a
 121 single, known behavior policy, we instead make the following regularity assumption:

122 **Assumption 1** (Reference distribution property). *For any (s, a) , $d^\pi(s, a) > 0$ implies $d^\mathcal{D}(s, a) > 0$.
 123 Furthermore, the correction terms are bounded by some finite constant C : $\|w_{\pi/\mathcal{D}}\|_\infty \leq C$.*

124 3 DualDICE

125 We now develop our algorithm, DualDICE, for estimating the discounted stationary distribution
 126 corrections $w_{\pi/\mathcal{D}}(s, a) = \frac{d^\pi(s, a)}{d^\mathcal{D}(s, a)}$. In the OPE setting, one does not have explicit knowledge of
 127 the distribution $d^\mathcal{D}$, but rather only access to samples $\mathcal{D} = \{(s, a, r, s')\} \sim d^\mathcal{D}$. Similar to the TD
 128 methods described above, we also assume access to samples from the initial state distribution β . We
 129 begin by introducing a key result, which we will later derive and use as the crux for our algorithm.

3.1 The Key Idea

Consider optimizing a (bounded) function $\nu : S \times A \rightarrow \mathbb{R}$ for the following objective:

$$\min_{\nu: S \times A \rightarrow \mathbb{R}} J(\nu) := \frac{1}{2} \mathbb{E}_{(s,a) \sim d^{\mathcal{D}}} [(\nu - \mathcal{B}^{\pi} \nu)(s, a)^2] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)], \quad (6)$$

where we use \mathcal{B}^{π} to denote the expected Bellman operator with respect to policy π and zero reward: $\mathcal{B}^{\pi} \nu(s, a) = \gamma \mathbb{E}_{s' \sim T(s, a), a' \sim \pi(s')} [\nu(s', a')]$. The first term in equation 6 is the expected squared Bellman error with zero reward. This term alone would lead to a trivial solution $\nu^* \equiv 0$, which can be avoided by the second term that encourages $\nu^* > 0$. Together, these two terms result in an optimal ν^* that places some non-zero amount of Bellman residual at state-action pairs sampled from $d^{\mathcal{D}}$.

Perhaps surprisingly, as we will show, the Bellman residuals of ν^* are exactly the desired distribution corrections:

$$(\nu^* - \mathcal{B}^{\pi} \nu^*)(s, a) = w_{\pi/\mathcal{D}}(s, a). \quad (7)$$

This key result provides the foundation for our algorithm, since it provides us with a simple objective (relying only on samples from $d^{\mathcal{D}}$, β , π) which we may optimize in order to obtain estimates of the distribution corrections. In the text below, we will show how we arrive at this result. We provide one additional step which allows us to efficiently learn a parameterized ν with respect to equation 6. We then generalize our results to a family of similar algorithms and lastly present theoretical guarantees.

3.2 Derivation

A Technical Observation We begin our derivation of the algorithm for estimating $w_{\pi/\mathcal{D}}$ by presenting the following simple technical observation: For arbitrary scalars $m \in \mathbb{R}_{>0}, n \in \mathbb{R}_{\geq 0}$, the optimizer of the convex problem $\min_x J(x) := \frac{1}{2} m x^2 - n x$ is unique and given by $x^* = \frac{n}{m}$. Using this observation, and letting \mathcal{C} be some bounded subset of \mathbb{R} which contains $[0, C]$, one immediately sees that the optimizer of the following problem,

$$\min_{x: S \times A \rightarrow \mathcal{C}} J_1(x) := \frac{1}{2} \mathbb{E}_{(s,a) \sim d^{\mathcal{D}}} [x(s, a)^2] - \mathbb{E}_{(s,a) \sim d^{\pi}} [x(s, a)], \quad (8)$$

is given by $x^*(s, a) = w_{\pi/\mathcal{D}}(s, a)$ for any $(s, a) \in S \times A$. This result provides us with an objective that shares the same basic form as equation 6. The main distinction is that the second term relies on an expectation over d^{π} , which we do not have access to.

Change of Variables In order to transform the second expectation in equation 8 to be over the initial state distribution β , we perform the following change of variables: Let $\nu : S \times A \rightarrow \mathbb{R}$ be an arbitrary state-action value function that satisfies,

$$\nu(s, a) := x(s, a) + \gamma \mathbb{E}_{s' \sim T(s, a), a' \sim \pi(s')} [\nu(s', a')], \quad \forall (s, a) \in S \times A. \quad (9)$$

Since $x(s, a) \in \mathcal{C}$ is bounded and $\gamma < 1$, the variable $\nu(s, a)$ is well-defined and bounded. By applying this change of variables, the objective function in 8 can be re-written in terms of ν , and this yields our previously presented objective from equation 6. Indeed, define,

$$\beta_t(s) := \Pr(s = s_t \mid s_0 \sim \beta, a_k \sim \pi(s_k), s_{k+1} \sim T(s_k, a_k) \text{ for } 0 \leq k < t),$$

to be the state visitation probability at step t when following π . Clearly, $\beta_0 = \beta$. Then,

$$\begin{aligned} \mathbb{E}_{(s,a) \sim d^{\pi}} [x(s, a)] &= \mathbb{E}_{(s,a) \sim d^{\pi}} [\nu(s, a) - \gamma \mathbb{E}_{s' \sim T(s, a), a' \sim \pi(s')} [\nu(s', a')]] \\ &= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim \beta_t, a \sim \pi(s)} [\nu(s, a) - \gamma \mathbb{E}_{s' \sim T(s, a), a' \sim \pi(s')} [\nu(s', a')]] \\ &= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim \beta_t, a \sim \pi(s)} [\nu(s, a)] - (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{s \sim \beta_{t+1}, a \sim \pi(s)} [\nu(s, a)] \\ &= (1 - \gamma) \mathbb{E}_{s \sim \beta, a \sim \pi(s)} [\nu(s, a)]. \end{aligned}$$

The Bellman residuals of the optimum of this objective give the desired off-policy corrections:

$$(\nu^* - \mathcal{B}^{\pi} \nu^*)(s, a) = x^*(s, a) = w_{\pi/\mathcal{D}}(s, a). \quad (10)$$

Equation 6 provides a promising approach for estimating the stationary distribution corrections, since the first expectation is over state-action pairs sampled from $d^{\mathcal{D}}$, while the second expectation is over β and actions sampled from π , both of which we have access to. Therefore, in principle we may solve this problem with respect to a parameterized value function ν , and then use the optimized ν^* to deduce the corrections. In practice, however, the objective in its current form presents two difficulties:

- The quantity $(\nu - \mathcal{B}^\pi \nu)(s, a)^2$ involves a conditional expectation inside a square. In general, when environment dynamics are stochastic and the action space may be large or continuous, this quantity may not be readily optimized using standard stochastic techniques. (For example, when the environment is stochastic, its Monte-Carlo sample gradient is generally biased.)
- Even if one has computed the optimal value ν^* , the corrections $(\nu^* - \mathcal{B}^\pi \nu^*)(s, a)$, due to the same argument as above, may not be easily computed, especially when the environment is stochastic or the action space continuous.

Exploiting Fenchel Duality We solve both difficulties listed above in one step by exploiting Fenchel duality [42]: Any convex function $f(x)$ may be written as $f(x) = \max_{\zeta} x \cdot \zeta - f^*(\zeta)$, where f^* is the Fenchel conjugate of f . In the case of $f(x) = \frac{1}{2}x^2$, the Fenchel conjugate is given by $f^*(\zeta) = \frac{1}{2}\zeta^2$. Thus, we may express our objective as,

$$\min_{\nu: S \times A \rightarrow \mathbb{R}} J(\nu) := \mathbb{E}_{(s,a) \sim d^D} \left[\max_{\zeta} (\nu - \mathcal{B}^\pi \nu)(s, a) \cdot \zeta - \frac{1}{2}\zeta^2 \right] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)].$$

By the interchangeability principle [8, 41, 43], we may replace the inner max over scalar ζ to a max over functions $\zeta: S \times A \rightarrow \mathbb{R}$ and obtain a min-max saddle-point optimization:

$$\min_{\nu: S \times A \rightarrow \mathbb{R}} \max_{\zeta: S \times A \rightarrow \mathbb{R}} J(\nu, \zeta) := \mathbb{E}_{(s,a,s') \sim d^D, a' \sim \pi(s')} [(\nu(s, a) - \gamma \nu(s', a'))\zeta(s, a) - \zeta(s, a)^2/2] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)]. \quad (11)$$

Using the KKT condition of the inner optimization problem (which is convex and quadratic in ζ), for any ν the optimal value ζ_ν^* is equal to the Bellman residual, $\nu - \mathcal{B}^\pi \nu$. Therefore, the desired stationary distribution correction can then be found from the saddle-point solution (ν^*, ζ^*) of the minimax problem in equation 11 as follows:

$$\zeta^*(s, a) = (\nu^* - \mathcal{B}^\pi \nu^*)(s, a) = w_{\pi/\mathcal{D}}(s, a). \quad (12)$$

Now we finally have an objective which is well-suited for practical computation. First, unbiased estimates of both the objective and its gradients are easy to compute using stochastic samples from d^D , β , and π , all of which we have access to. Secondly, notice that the min-max objective function in equation 11 is linear in ν and concave in ζ . Therefore in certain settings, one can provide guarantees on the convergence of optimization algorithms applied to this objective (see Section 3.4). Thirdly, the optimizer of the objective in equation 11 immediately gives us the desired stationary distribution corrections through the values of $\zeta^*(s, a)$, with no additional computation.

3.3 Extension to General Convex Functions

Besides a quadratic penalty function, one may extend the above derivations to a more general class of convex penalty functions. Consider a generic convex penalty function $f: \mathbb{R} \rightarrow \mathbb{R}$. Recall that \mathcal{C} is a bounded subset of \mathbb{R} which contains the interval $[0, C]$ of stationary distribution corrections. If \mathcal{C} is contained in the range of f' , then the optimizer of the convex problem, $\min_x J(x) := m \cdot f(x) - n$ for $\frac{n}{m} \in \mathcal{C}$, satisfies the following KKT condition: $f'(x^*) = \frac{n}{m}$. Analogously, the optimizer x^* of,

$$\min_{x: S \times A \rightarrow \mathcal{C}} J_1(x) := \mathbb{E}_{(s,a) \sim d^D} [f(x(s, a))] - \mathbb{E}_{(s,a) \sim d^\pi} [x(s, a)], \quad (13)$$

satisfies the equality $f'(x^*(s, a)) = w_{\pi/\mathcal{D}}(s, a)$.

With change of variables $\nu := x + \mathcal{B}^\pi \nu$, the above problem becomes,

$$\min_{\nu: S \times A \rightarrow \mathbb{R}} J(\nu) := \mathbb{E}_{(s,a) \sim d^D} [f((\nu - \mathcal{B}^\pi \nu)(s, a))] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)]. \quad (14)$$

Applying Fenchel duality to f in this objective further leads to the following saddle-point problem:

$$\min_{\nu: S \times A \rightarrow \mathbb{R}} \max_{\zeta: S \times A \rightarrow \mathbb{R}} J(\nu, \zeta) := \mathbb{E}_{(s,a,s') \sim d^D, a' \sim \pi(s')} [(\nu(s, a) - \gamma \nu(s', a'))\zeta(s, a) - f^*(\zeta(s, a))] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)]. \quad (15)$$

By the KKT condition of the inner optimization problem, for any ν the optimizer ζ_ν^* satisfies,

$$f^{*'}(\zeta_\nu^*(s, a)) = (\nu - \mathcal{B}^\pi \nu)(s, a). \quad (16)$$

Therefore, using the fact that the derivative of a convex function f' is the inverse function of the derivative of its Fenchel conjugate $f^{*'}$, our desired stationary distribution corrections are found by computing the saddle-point (ζ^*, ν^*) of the above problem:

$$\zeta^*(s, a) = f'((\nu^* - \mathcal{B}^\pi \nu^*)(s, a)) = f'(x^*(s, a)) = w_{\pi/\mathcal{D}}(s, a). \quad (17)$$

Amazingly, despite the generalization beyond the quadratic penalty function $f(x) = \frac{1}{2}x^2$, the optimization problem in equation 15 retains all the computational benefits that make this method very practical for learning $w_{\pi/\mathcal{D}}(s, a)$: All quantities and their gradients may be unbiasedly estimated from stochastic samples; the objective is linear in ν and concave in ζ , thus is well-behaved; and the optimizer of this problem immediately provides the desired stationary distribution corrections through the values of $\zeta^*(s, a)$, without any additional computation.

This generalized derivation also provides insight into the initial technical result: It is now clear that the objective in equation 13 is the negative Fenchel dual (variational) form of the Ali-Silvey or f -divergence, which has been used in previous work to estimate divergence and data likelihood ratios [33]. Despite their similar formulations, we emphasize that the aforementioned dual form of the f -divergence is not immediately applicable to estimation of off-policy corrections in the context of RL, due to the fact that samples from distribution d^π are unobserved. Indeed, our derivations hinged on two additional key steps: (1) the change of variables from x to $\nu := x + \mathcal{B}^\pi \nu$; and (2) the second application of duality to introduce ζ . Due to these repeated applications of duality in our derivations, we term our method *Dual stationary DIstribution Correction Estimation (DualDICE)*.

3.4 Theoretical Guarantees

In this section, we consider the theoretical properties of DualDICE in the setting where we have a dataset formed by empirical samples $\{s_i, a_i, r_i, s'_i\}_{i=1}^N \sim d^\mathcal{D}$, $\{s_0^i\}_{i=1}^N \sim \beta$, and target actions $a'_i \sim \pi(s'_i)$, $a_0^i \sim \pi(s_0^i)$ for $i = 1, \dots, N$.¹ We will use the shorthand notation $\hat{\mathbb{E}}_{d^\mathcal{D}}$ to denote an average over these empirical samples. Although the proposed estimator can adopt general f , for simplicity of exposition we restrict to $f(x) = \frac{1}{2}x^2$. We consider using an algorithm *OPT* (e.g., stochastic gradient descent/ascent) to find optimal ν, ζ of equation 15 within some parameterization families \mathcal{F}, \mathcal{H} , respectively. We denote by $\hat{\nu}, \hat{\zeta}$ the outputs of *OPT*. We have the following guarantee on the quality of $\hat{\nu}, \hat{\zeta}$ with respect to the off-policy policy estimation (OPE) problem.

Theorem 2. (Informal) Under some mild assumptions, the mean squared error (MSE) associated with using $\hat{\nu}, \hat{\zeta}$ for OPE can be bounded as,

$$\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^\mathcal{D}} \left[\hat{\zeta}(s, a) \cdot r \right] - \rho(\pi) \right)^2 \right] = \tilde{\mathcal{O}} \left(\epsilon_{\text{approx}}(\mathcal{F}, \mathcal{H}) + \epsilon_{\text{opt}} + \frac{1}{\sqrt{N}} \right), \quad (18)$$

where the outer expectation is with respect to the randomness of the empirical samples and *OPT*, ϵ_{opt} denotes the optimization error, and $\epsilon_{\text{approx}}(\mathcal{F}, \mathcal{H})$ denotes the approximation error due to \mathcal{F}, \mathcal{H} .

The sources of estimation error are explicit in Theorem 2. As the number of samples N increases, the statistical error $N^{-1/2}$ approaches zero. Meanwhile, there is an implicit trade-off in $\epsilon_{\text{approx}}(\mathcal{F}, \mathcal{H})$ and ϵ_{opt} . With flexible function spaces \mathcal{F} and \mathcal{H} (such as the space of neural networks), the approximation error can be further decreased; however, optimization will be complicated and it is difficult to characterize ϵ_{opt} . On the other hand, with linear parameterization of (ν, ζ) , under some mild conditions, after T iterations we achieve provably fast rate, $\mathcal{O}(\exp(-T))$ for *OPT* = SVRG and $\mathcal{O}(\frac{1}{T})$ for *OPT* = SGD, at the cost of potentially increased approximation error. See the Appendix for the precise theoretical results, proofs, and further discussions.

4 Related Work

Density Ratio Estimation Density ratio estimation is an important tool for many machine learning and statistics problems. Other than the naive approach, (i.e., the density ratio is calculated via estimating the densities in the numerator and denominator separately, which may magnify the estimation error), various direct ratio estimators have been proposed [44], including the moment matching approach [17], probabilistic classification approach [4, 7, 40], and ratio matching approach [22, 33, 45]

The proposed DualDICE algorithm, as a direct approach for density ratio estimation, bears some similarities to ratio matching [33], which is also derived by exploiting the Fenchel dual representation of the f -divergences. However, compared to the existing direct estimators, the major difference lies in the requirement of the samples from the stationary distribution. Specifically, the existing estimators require access to samples from both $d^\mathcal{D}$ and d^π , which is impractical in the off-policy learning setting. Therefore, DualDICE is uniquely applicable to the more difficult RL setting.

¹For the sake of simplicity, we consider the batch learning setting with *i.i.d.* samples as in [48]. The results can be easily generalized to single sample path with dependent samples (see Appendix).

Off-policy Policy Evaluation The problem of off-policy policy evaluation has been heavily studied in contextual bandits [12, 49, 52] and in the more general RL setting [14, 21, 26, 29, 34, 36, 37, 50, 51]. Several representative approaches can be identified in the literature. The Direct Method (DM) learns a model of the system and then uses it to estimate the performance of the evaluation policy. This approach often has low variance but its bias depends on how well the selected function class can express the environment dynamics. Importance sampling (IS) [38] uses importance weights to correct the mismatch between the distributions of the system trajectory induced by the target and behavior policies. Its variance can be unbounded when there is a big difference between the distributions of the evaluation and behavior policies, and grows exponentially with the horizon of the RL problem. Doubly Robust (DR) is a combination of DM and IS, and can achieve the low variance of DM and no (or low) bias of IS. Other than DM, all the methods described above require knowledge of the policy density ratio, and thus the behavior policy. Our proposed algorithm avoids this necessity.

5 Experiments

We evaluate our method applied to off-policy policy evaluation (OPE). We focus on this setting because it is a direct application of stationary distribution correction estimation, without many additional tunable parameters, and it has been previously used as a test-bed for similar techniques [27]. In each experiment, we use a behavior policy μ to collect some number of trajectories, each for some number of steps. This data is used to estimate the stationary distribution corrections, which are then used to estimate the average step reward, with respect to a target policy π . We focus our comparisons here to a TD-based approach [16] and weighted step-wise IS (as described in [27]), which we and others have generally found to work best relative to common IS variants [30, 38]. See the Appendix for additional results and implementation details.

We begin in a controlled setting with an evaluation agnostic to optimization issues, where we find that, absent these issues, our method is competitive with TD-based approaches (Figure 1). However, as we move to more difficult settings with complex environment dynamics, the performance of TD methods degrades dramatically, while our method is still able to provide accurate estimates (Figure 2). Finally, we provide an analysis of the optimization behavior of our method on a simple control task across different choices of function f (Figure 3). Interestingly, although the choice of $f(x) = \frac{1}{2}x^2$ is most natural, we find the empirically best performing choice to be $f(x) = \frac{2}{3}|x|^{3/2}$. All results are summarized for 20 random seeds, with median plotted and error bars at 25th and 75th percentiles.

5.1 Estimation Without Function Approximation

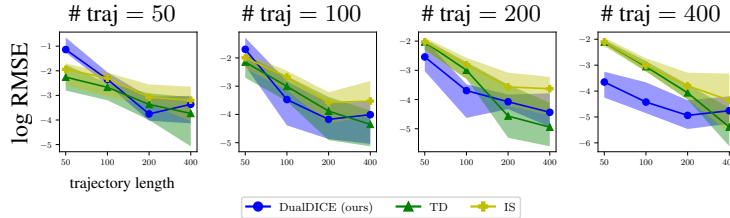


Figure 1: We perform OPE on the Taxi domain [10]. The plots show log RMSE of the estimator across different numbers of trajectories and different trajectory lengths (x -axis). For this domain, we avoid any potential issues in optimization by solving for the optimum of the objectives exactly using standard matrix operations. Thus, we are able to see that our method and the TD method are competitive with each other.

We begin with a tabular task, the Taxi domain [10]. In this task, we evaluate our method in a manner agnostic to optimization difficulties: The objective J is a quadratic equation in ν , and thus may be solved by matrix operations. The Bellman residuals (equation 7) may then be estimated via an empirical average of the transitions appearing in the off-policy data. In a similar manner, TD methods for estimating the correction terms may also be solved using matrix operations [27]. In this controlled setting, we find that, as expected, TD methods can perform well (Figure 1), and our method achieves competitive performance. As we will see in the following results, the good performance of TD methods quickly deteriorates as one moves to more complex settings, while our method is able to maintain good performance, even when using function approximation and stochastic optimization.

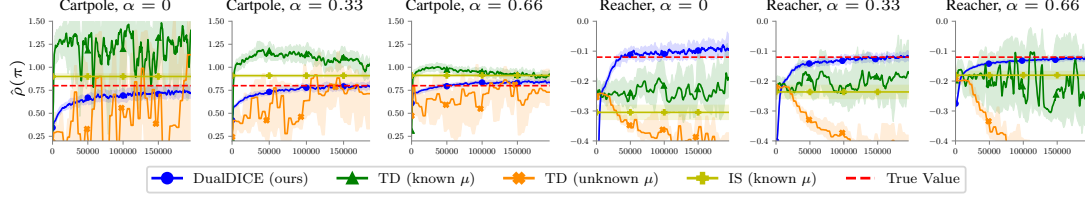


Figure 2: We perform OPE on control tasks. Each plot shows the estimated average step reward over training and different behavior policies (higher α corresponds to a behavior policy closer to the target policy). We find that in all cases, our method is able to approximate these desired values well, with accuracy improving with a larger α . On the other hand, the TD method performs poorly, even more so when the behavior policy μ is unknown and must be estimated. While on Cartpole it can start to approach the desired value for large α , on the more complicated Reacher task (which involves continuous actions) its learning is too unstable to learn anything at all.

5.2 Control Tasks

We now move on to difficult control tasks: A discrete-control task Cartpole and a continuous-control task Reacher [6]. In these tasks, observations are continuous, and thus we use neural network function approximators with stochastic optimization. Figure 2 shows the results of our method compared to the TD method. We find that in this setting, DualDICE is able to provide good, stable performance, while the TD approach suffers from high variance, and this issue is exacerbated when we attempt to estimate μ rather than assume it as given. See the Appendix for additional baseline results.

5.3 Choice of Convex Function f

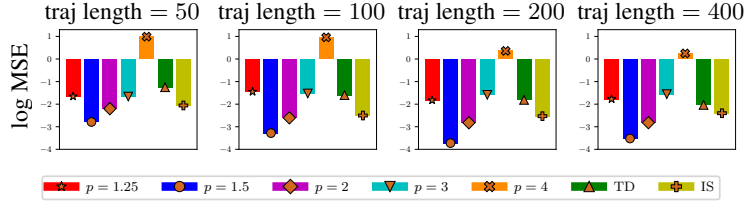


Figure 3: We compare the OPE error when using different forms of f to estimate stationary distribution ratios with function approximation, which are then applied to OPE on a simple continuous grid task. In this setting, optimization stability is crucial, and this heavily depends on the form of the convex function f . We plot the results of using $f(x) = \frac{1}{p}|x|^p$ for $p \in [1.25, 1.5, 2, 3, 4]$. We also show the results of TD and IS methods on this task for comparison. We find that $p = 1.5$ consistently performs the best, often providing significantly better results.

We analyze the choice of the convex function f . We consider a simple continuous grid task where an agent may move left, right, up, or down and is rewarded for reaching the bottom right corner of a square room. We plot the estimation errors of using DualDICE for off-policy policy evaluation on this task, comparing against different choices of convex functions of the form $f(x) = \frac{1}{p}|x|^p$. Interestingly, although the choice of $f(x) = \frac{1}{2}x^2$ is most natural, we find the empirically best performing choice to be $f(x) = \frac{2}{3}|x|^{3/2}$. Thus, this is the form of f we used in our experiments for Figure 2.

6 Conclusions

We have presented DualDICE, a method for estimating off-policy stationary distribution corrections. Compared to previous work, our method is agnostic to knowledge of the behavior policy used to collect the off-policy data and avoids the use of importance weights in its losses. These advantages have a profound empirical effect: our method provides significantly better estimates compared to TD methods, especially in settings which require function approximation and stochastic optimization.

Future work includes (1) incorporating the DualDICE algorithm into off-policy training, (2) further understanding the effects of f on the performance of DualDICE (in terms of approximation error of the distribution corrections), and (3) evaluating DualDICE on real-world off-policy evaluation tasks.

References

- [1] Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018.
- [2] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- [3] Richard Ernest Bellman. *Dynamic Programming*. Dover Publications, Inc., New York, NY, USA, 2003.
- [4] Steffen Bickel, Michael Brückner, and Tobias Scheffer. Discriminative learning for differing training and test distributions. In *Proceedings of the 24th international conference on Machine learning*, pages 81–88. ACM, 2007.
- [5] Stephane Boucheron, Gabor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2016.
- [6] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [7] Kuang Fu Cheng, Chih-Kang Chu, et al. Semiparametric density estimation under a two-sample density ratio model. *Bernoulli*, 10(4):583–604, 2004.
- [8] Bo Dai, Niao He, Yunpeng Pan, Byron Boots, and Le Song. Learning from conditional distributions via dual embeddings. *arXiv preprint arXiv:1607.04579*, 2016.
- [9] Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbed: Convergent reinforcement learning with nonlinear function approximation. *arXiv preprint arXiv:1712.10285*, 2017.
- [10] Thomas G Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- [11] Simon S Du, Jianshu Chen, Lihong Li, Lin Xiao, and Dengyong Zhou. Stochastic variance reduction methods for policy evaluation. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1049–1058. JMLR. org, 2017.
- [12] Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. *arXiv preprint arXiv:1103.4601*, 2011.
- [13] Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly robust off-policy evaluation. *arXiv preprint arXiv:1802.03493*, 2018.
- [14] Raphael Fonteneau, Susan A. Murphy, Louis Wehenkel, and Damien Ernst. Batch mode reinforcement learning based on the synthesis of artificial trajectories. *Annals of Operations Research*, 208(1):383–416, 2013.
- [15] Jianfeng Gao, Michel Galley, and Lihong Li. Neural approaches to Conversational AI. *Foundations and Trends in Information Retrieval*, 13(2–3):127–298, 2019.
- [16] Carles Gelada and Marc G Bellemare. Off-policy deep reinforcement learning by bootstrapping the covariate shift. *AAAI*, 2018.
- [17] Arthur Gretton, Alex J Smola, Jiayuan Huang, Marcel Schmittfull, Karsten M Borgwardt, and Bernhard Schölkopf. Covariate shift by kernel mean matching. In *Dataset shift in machine learning*, pages 131–160. MIT Press, 2009.
- [18] Assaf Hallak and Shie Mannor. Consistent on-line off-policy evaluation. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1372–1383. JMLR. org, 2017.

- [19] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications. 1970.
- [20] David Haussler. Sphere packing numbers for subsets of the boolean n-cube with bounded vapnik-chervonenkis dimension. *Journal of Combinatorial Theory, Series A*, 69(2):217–232, 1995.
- [21] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 652–661, 2016.
- [22] Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10(Jul):1391–1445, 2009.
- [23] Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of least-squares policy iteration. *Journal of Machine Learning Research*, 13(Oct):3041–3074, 2012.
- [24] Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, and Dan Jurafsky. Deep reinforcement learning for dialogue generation. *arXiv preprint arXiv:1606.01541*, 2016.
- [25] Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-bandit-based news article recommendation algorithms. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 297–306. ACM, 2011.
- [26] Lihong Li, Rémi Munos, and Csaba Szepesvári. Toward minimax off-policy value estimation. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, pages 608–616, 2015.
- [27] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pages 5356–5366, 2018.
- [28] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient with state distribution correction. In *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence*, 2019. To appear.
- [29] A. Mahmood, H. van Hasselt, and R. Sutton. Weighted importance sampling for off-policy learning with linear function approximation. In *Proceedings of the 27th International Conference on Neural Information Processing Systems*, 2014.
- [30] Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pages 1077–1084. International Foundation for Autonomous Agents and Multiagent Systems, 2014.
- [31] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [32] Susan A Murphy, Mark J van der Laan, James M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423, 2001.
- [33] XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- [34] C. Paduraru. *Off-policy Evaluation in Markov Decision Processes*. PhD thesis, McGill University, 2013.
- [35] D Pollard. *Convergence of Stochastic Processes*. David Pollard, 1984.

- 402 [36] D. Precup, R. Sutton, and S. Dasgupta. Off-policy temporal difference learning with function
403 approximation. In *Proceedings of the 18th International Conference on Machine Learning*,
404 pages 417–424, 2001.
- 405 [37] D. Precup, R. Sutton, and S. Singh. Eligibility traces for off-policy policy evaluation. In
406 *Proceedings of the 17th International Conference on Machine Learning*, pages 759–766, 2000.
- 407 [38] Doina Precup. Eligibility traces for off-policy policy evaluation. *Computer Science Department*
408 *Faculty Publication Series*, page 80, 2000.
- 409 [39] Martin L Puterman. Markov decision processes: Discrete stochastic dynamic programming.
410 1994.
- 411 [40] Jing Qin. Inferences for case-control and semiparametric two-sample density ratio models.
412 *Biometrika*, 85(3):619–630, 1998.
- 413 [41] R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science
414 & Business Media, 2009.
- 415 [42] Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015.
- 416 [43] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic*
417 *programming: modeling and theory*. SIAM, 2009.
- 418 [44] Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density ratio estimation in machine*
419 *learning*. Cambridge University Press, 2012.
- 420 [45] Masashi Sugiyama, Taiji Suzuki, Shinichi Nakajima, Hisashi Kashima, Paul von Bünau, and
421 Motoaki Kawanabe. Direct importance estimation for covariate shift adaptation. *Annals of the*
422 *Institute of Statistical Mathematics*, 60(4):699–746, 2008.
- 423 [46] Richard S Sutton and Andrew G Barto. *Introduction to reinforcement learning*, volume 135.
- 424 [47] Richard S Sutton, A Rupam Mahmood, and Martha White. An emphatic approach to the
425 problem of off-policy temporal-difference learning. *The Journal of Machine Learning Research*,
426 17(1):2603–2631, 2016.
- 427 [48] Richard S Sutton, Csaba Szepesvári, Alborz Geramifard, and Michael Bowling. Dyna-style
428 planning with linear function approximation and prioritized sweeping. In *Proceedings of the*
429 *Twenty-Fourth Conference on Uncertainty in Artificial Intelligence*, pages 528–536. AUAI Press,
430 2008.
- 431 [49] A. Swaminathan, A. Krishnamurthy, A. Agarwal, M. Dudík, J. Langford, D. Jose, and I. Zitouni.
432 Off-policy evaluation for slate recommendation. In *Proceedings of the 31st International*
433 *Conference on Neural Information Processing Systems*, pages 3635–3645, 2017.
- 434 [50] P. Thomas and E. Brunskill. Data-efficient off-policy policy evaluation for reinforcement
435 learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages
436 2139–2148, 2016.
- 437 [51] P. Thomas, G. Theodorou, and M. Ghavamzadeh. High confidence off-policy evaluation. In
438 *Proceedings of the 29th Conference on Artificial Intelligence*, 2015.
- 439 [52] Yu-Xiang Wang, Alekh Agarwal, and Miroslav Dudík. Optimal and adaptive off-policy evalua-
440 tion in contextual bandits. In *Proceedings of the 34th International Conference on Machine*
441 *Learning-Volume 70*, pages 3589–3597. JMLR. org, 2017.
- 442 [53] Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The*
443 *Annals of Probability*, pages 94–116, 1994.

444 A Pseudocode

Algorithm 1 DualDICE

Inputs: Convex function f and its Fenchel conjugate f^* , off-policy data $\hat{\mathcal{D}} = \{(s^{(i)}, a^{(i)}, r^{(i)}, s'^{(i)})\}_{i=1}^N$, sampled initial states $\hat{\beta} = \{s_0^{(i)}\}_{i=1}^M$, target policy π , networks $\nu_{\theta_1}(\cdot, \cdot), \zeta_{\theta_2}(\cdot, \cdot)$, learning rates η_ν, η_ζ , number of iterations T , batch size B .

for $t = 1, \dots, T$ **do**

Sample batch $\{(s^{(i)}, a^{(i)}, r^{(i)}, s'^{(i)})\}_{i=1}^B$ from $\hat{\mathcal{D}}$.

Sample batch $\{s_0^{(i)}\}_{i=1}^B$ from $\hat{\beta}$.

Sample actions $a'^{(i)} \sim \pi(s'^{(i)})$, for $i = 1, \dots, B$.

Sample actions $a_0^{(i)} \sim \pi(s_0^{(i)})$, for $i = 1, \dots, B$.

Compute empirical loss $\hat{J} = \frac{1}{B} \sum_{i=1}^B (\nu_{\theta_1}(s^{(i)}, a^{(i)}) - \nu_{\theta_1}(s'^{(i)}, a'^{(i)})) \zeta_{\theta_2}(s^{(i)}, a^{(i)}) - f^*(\zeta_{\theta_2}(s^{(i)}, a^{(i)})) - (1 - \gamma) \nu_{\theta_1}(s_0^{(i)}, a_0^{(i)})$.

Update $\theta_1 \leftarrow \theta_1 - \eta_\nu \nabla_{\theta_1} \hat{J}$.

Update $\theta_2 \leftarrow \theta_2 + \eta_\zeta \nabla_{\theta_2} \hat{J}$.

end for

Return $\zeta_{\theta_2}(\cdot, \cdot)$.

445 B Additional Results

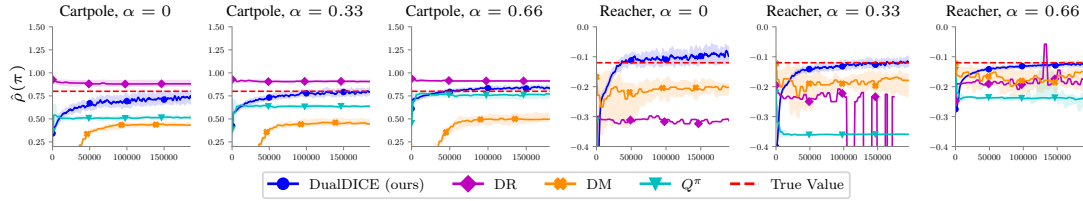


Figure 4: We perform OPE on control tasks using our method compared to a number of additional baselines: doubly-robust (DR), in which one learns a value function in order to reduce the variance of an IS estimate of the evaluation; direct method (DM), in which one learns a model of the dynamics and reward of the environment and performs Monte Carlo rollouts using the model in order to estimate the value of the target policy; and Q^π , in which one learns Q^π values via Bellman error minimization over the off-policy data, and uses the initial values $(1 - \gamma) \cdot Q^\pi(s_0, a_0)$ as estimates of the policy value (these estimates are below -0.4 for Reacher, $\alpha = 0$).

446 C Experimental Details

447 C.1 Taxi

448 For the Taxi domain, we follow the same protocol as used in [27]. In this tabular, exact solve setting,
 449 the TD methods [16] are equivalent to their kernel-based TD method. We fix γ to 0.995. The behavior
 450 and target policies are also taken from [27] (referred in their work as the behavior policy for $\alpha = 0$).

451 In this setting, we solve for the optimal empirical ν exactly using matrix operations. Since [27]
 452 perform a similar exact solve for $|S|$ variables $w_{\pi/\mu}(s)$, for better comparison we also perform our
 453 exact solve with respect to $|S|$ variables $\nu(s)$. Specifically, one may follow the same derivations
 454 for DualDICE with respect to learning $w_{\pi/\mu}$. The final objective will require knowledge of the
 455 importance weights $\pi(a|s)/\mu(a|s)$.

456 C.2 Control Tasks

457 We use the Cartpole and Reacher tasks as given by OpenAI Gym [6]. In these tasks we use COP-
 458 TD [16] for the TD method ([27] requires a proper kernel, which is not readily available for these

tasks). When assuming an unknown μ , we learn a neural network policy $\hat{\mu}$ using behavior cloning, and use its probabilities for computing importance weights $\pi(a|s)/\mu(a|s)$. All neural networks are feed-forward with two hidden layers of dimension 64 and tanh activations.

We modify the Cartpole task to be infinite horizon: We use the same dynamics as in the original task but change the reward to be -1 if the original task returns a termination (when the pole falls below some threshold) and 1 otherwise. We train a policy on this task until convergence. We then define the target policy π as a weighted combination of this pre-trained policy (weight 0.7) and a uniformly random policy (weight 0.3). The behavior policy μ for a specific $0 \leq \alpha \leq 1$ is taken to be a weighted combination of the pre-trained policy (weight $0.55 + 0.15\alpha$) and a uniformly random policy (weight $0.45 - 0.15\alpha$). We use $\gamma = 0.99$, which yields an average step reward of ≈ 0.8 for π and ≈ 0.1 for μ with $\alpha = 0$. We generate an off-policy dataset by running the behavior policy for 200 episodes, each of length 250 steps. We train each stationary distribution correction estimation method using the Adam optimizer with batches of size 2048 and learning rates chosen using a hyperparameter search (the optimal learning rate found for either method was ≈ 0.003).

For the Reacher task, we train a deterministic policy until convergence. We define the target policy π as a Gaussian with mean given by the pre-trained policy and standard deviation given by 0.1 . The behavior policy μ for a specific $0 \leq \alpha \leq 1$ is taken to be a Gaussian with mean given by the pre-trained policy and standard deviation given by $0.4 - 0.3\alpha$. We use $\gamma = 0.99$, which yields an average step reward of ≈ -0.12 for π and ≈ -0.50 for μ with $\alpha = 0$. We generate an off-policy dataset by running the behavior policy for 1000 episodes, each of length 40 steps. We train each stationary distribution correction estimation method using the Adam optimizer with batches of size 2048 and learning rates chosen using a hyperparameter search (the optimal learning rate found for either method was ≈ 0.0001).

C.3 Continuous Grid

For this task, we create a 10×10 grid which the agent can traverse by moving left/right/up/down. The observations are the x, y coordinates of the square the agent is on. The reward at each step is given by $\exp\{-0.2|x - 9| - 0.2|y - 9|\}$. We use $\gamma = 0.995$. The target policy π is taken to be the optimal policy for this task plus 0.1 weight on uniform exploration. The behavior policy μ is taken to be the optimal policy plus 0.7 weight on uniform exploration. We train using batches of size the Adam optimizer with batches of size 512 and learning rates 0.001 for ν and 0.0001 for ζ .

D Proofs

We provide the proof for Theorem 2. We first decompose the error in Section D.1. Then, we analyze the statistical error and optimization error in Section D.2 and Section D.4, respectively. The total error will be discussed in D.3.

Although the proposed estimator can use any general convex function f , as a first step towards a more complete theoretical understanding, we consider the special case of $f(x) = \frac{1}{2}x^2$. Clearly, $f(\cdot)$ now is η -strongly convex with $\eta = 1$. Under Assumption 1, we need only consider $\|\nu\|_\infty \leq C$, which implies that $\|\nu - \mathcal{B}^\pi \nu\|_\infty \leq \frac{1+\gamma}{1-\gamma}C$, and that $f(x)$ is κ -Lipschitz continuous with $\kappa = \frac{1+\gamma}{1-\gamma}C$. Similarly, $f^*(y) = \frac{1}{2}y^2$ is L -Lipschitz continuous with $L = C$ on $\|w\|_\infty \leq C$. The following assumption will be needed.

Assumption 3 (MDP regularity). *We assume the observed reward is uniformly bounded, i.e., $\|\hat{r}(s, a)\|_\infty \leq C_r$ for some constant $C_r > 0$. It follows that the reward’s mean and variance are both bounded in $[-C_r, C_r]$.*

For convenience, the objective function of DualDICE is repeated here:

$$J(\nu, \zeta) = \mathbb{E}_{(s, a, s') \sim d^{\mathcal{D}}, a' \sim \pi(s')} [(\nu(s, a) - \gamma \nu(s', a'))(\zeta(s, a) - \zeta(s, a')^2/2] \\ - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)].$$

We will also make use of the objective in the form prior to introduction of ζ , which we denote as $J(\nu)$:

$$J(\nu) = \frac{1}{2} \mathbb{E}_{(s, a) \sim d^{\mathcal{D}}} [(\nu - \mathcal{B}^\pi \nu)(s, a)^2] - (1 - \gamma) \mathbb{E}_{s_0 \sim \beta, a_0 \sim \pi(s_0)} [\nu(s_0, a_0)].$$

Let $\hat{J}(\nu, \zeta)$ denotes the empirical surrogate of $J(\nu, \zeta)$ with optimal solution as $(\hat{\nu}^*, \hat{\zeta}^*)$. We denote $\nu_{\mathcal{F}}^* = \arg \min_{\nu \in \mathcal{F}} J(\nu)$ and $\nu^* = \arg \min_{\nu \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} J(\nu)$. We denote $L(\nu) = \max_{\zeta \in \mathcal{H}} J(\nu, \zeta)$ and $\hat{L}(\nu) = \max_{\zeta \in \mathcal{H}} \hat{J}(\nu, \zeta)$ as the primal objectives, and $\ell(\zeta) = \min_{\nu \in \mathcal{F}} J(\nu, \zeta)$, $\hat{\ell}(\zeta) = \min_{\nu \in \mathcal{F}} \hat{J}(\nu, \zeta)$ as the dual objectives. We apply some optimization algorithm *OPT* for optimizing $\hat{J}(\nu, \zeta)$ with samples $\{s_i, a_i, r_i, s'_i\}_{i=1}^N \sim d^{\mathcal{D}}$, $\{s_0^i\}_{i=1}^N \sim \beta$, and target actions $a'_i \sim \pi(s'_i)$, $a_0^i \sim \pi(s_0^i)$ for $i = 1, \dots, N$. We denote the outputs of *OPT* by $(\hat{\nu}, \hat{\zeta})$.

D.1 Error Decomposition

Let

$$\bar{R}(s, a) = \mathbb{E}_{\cdot|s,a} [r].$$

Observe that

$$\rho(\pi) = \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)].$$

We begin by considering the estimation error induced by using $(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a)$ as estimates of $w_{\pi/\mathcal{D}}(s, a)$, where $\hat{\mathcal{B}}^{\pi}$ denotes the empirical Bellman backup with respect to samples from $d^{\mathcal{D}}$, π . We will subsequently reconcile this with the true implementation of DualDICE, which uses $\hat{\zeta}(s, a)$ as estimates of $w_{\pi/\mathcal{D}}(s, a)$.

The mean squared error of the policy value estimate when using $(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a)$ in place of $w_{\pi/\mathcal{D}}(s, a)$ can be decomposed as

$$\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot r \right] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)] \right)^2 \quad (19)$$

$$= \left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot r \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \bar{R}(s, a) \right] \right. \\ \left. + \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \bar{R}(s, a) \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)(s, a) \cdot \bar{R}(s, a) \right] \right. \\ \left. + \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)(s, a) \cdot \bar{R}(s, a) \right] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)] \right)^2 \quad (20)$$

$$\leq 4 \underbrace{\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot r \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \bar{R}(s, a) \right] \right)^2}_{\epsilon_r} \quad (21)$$

$$+ 4 \underbrace{\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \bar{R}(s, a) \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)(s, a) \cdot \bar{R}(s, a) \right] \right)^2}_{\epsilon_1} \quad (22)$$

$$+ 4 \underbrace{\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)(s, a) \cdot \bar{R}(s, a) \right] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)] \right)^2}_{\epsilon_2}. \quad (23)$$

The first term, ϵ_r , is induced by the randomness in observed reward, and we have

$$\epsilon_r \leq \left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot (\hat{r}(s, a) - r(s, a)) \right] \right)^2 \leq \left(\frac{1+\gamma}{1-\gamma} \right)^2 C^2 \left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [\hat{r}(s, a)] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} [r(s, a)] \right)^2,$$

which will be discussed in section D.2.

We consider the ϵ_1 as

$$\epsilon_1 \leq C_r^2 \left\| (\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu}) - (\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*) \right\|_{\hat{\mathcal{D}}}^2 \leq C_r^2 \underbrace{\left(\left\| \hat{\zeta} - \hat{\zeta}^* \right\|_{\hat{\mathcal{D}}}^2 + \left\| (\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*) - (\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu}) \right\|_{\hat{\mathcal{D}}}^2 \right)}_{\hat{\epsilon}_{opt}}$$

which is the error induced by optimization *OPT*.

524 For the last term ϵ_2 , we have

$$\begin{aligned}
\epsilon_2 &\leq 2 \underbrace{\left(\hat{\mathbb{E}}_{d^D} \left[\left(\hat{\nu}^* - \hat{\mathcal{B}}^\pi \hat{\nu}^* \right) (s, a) \cdot r(s, a) \right] - \mathbb{E}_{d^D} \left[\left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) (s, a) \cdot r(s, a) \right] \right)^2}_{\epsilon_{stat}} \\
&\quad + 2 \left(\mathbb{E}_{d^D} \left[\left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) (s, a) \cdot r(s, a) \right] - \mathbb{E}_{d^D} \left[w_{\pi/D}(s, a) \cdot r(s, a) \right] \right)^2 \\
&\leq 2\epsilon_{stat} + 2 \left(\mathbb{E}_{d^D} \left[\left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) (s, a) \cdot r(s, a) \right] - \mathbb{E}_{d^D} \left[\left(\nu^* - \mathcal{B}^\pi \nu^* \right) (s, a) \cdot r(s, a) \right] \right)^2. \quad (\text{due to equation 17})
\end{aligned}$$

525 For the first term ϵ_{stat} , which is due to finite samples, we will bound in section D.2.

526 For the second term, we have

$$\begin{aligned}
&\left(\mathbb{E}_{d^D} \left[\left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) (s, a) \cdot r(s, a) \right] - \mathbb{E}_{d^D} \left[\left(\nu^* - \mathcal{B}^\pi \nu^* \right) (s, a) \cdot r(s, a) \right] \right)^2 \\
&\leq \mathbb{E}_{d^D} \left[r(s, a)^2 \cdot \left(\left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) (s, a) - \left(\nu^* - \mathcal{B}^\pi \nu^* \right) (s, a) \right)^2 \right] \\
&\leq C_r^2 \left\| \left(\hat{\nu}^* - \mathcal{B}^\pi \hat{\nu}^* \right) - \left(\nu^* - \mathcal{B}^\pi \nu^* \right) \right\|_{\mathcal{D}}^2 \\
&\leq \frac{2C_r^2}{\eta} \left(J(\hat{\nu}^*) - J(\nu^*) \right),
\end{aligned}$$

527 where the last inequality comes from the η -strongly convexity of f and the optimality of ν^* .

528 We then consider the error between $J(\hat{\nu}^*)$ and $J(\nu^*)$, which can be decomposed as

$$\begin{aligned}
J(\hat{\nu}^*) - J(\nu^*) &= J(\hat{\nu}^*) - J(\nu_{\mathcal{F}}^*) + J(\nu_{\mathcal{F}}^*) - J(\nu^*) \\
&= J(\hat{\nu}^*) - L(\hat{\nu}^*) + L(\hat{\nu}^*) - L(\nu_{\mathcal{F}}^*) + L(\nu_{\mathcal{F}}^*) - J(\nu_{\mathcal{F}}^*) + J(\nu_{\mathcal{F}}^*) - J(\nu^*).
\end{aligned}$$

529 We bound this expression term-by-term from the right. For the term $J(\nu_{\mathcal{F}}^*) - J(\nu^*)$, we have

$$\begin{aligned}
J(\nu_{\mathcal{F}}^*) - J(\nu^*) &= \mathbb{E}_{\mathcal{D}} [f(\nu_{\mathcal{F}}^* - \mathcal{B}^\pi \nu_{\mathcal{F}}^*) - f(\nu^* - \mathcal{B}^\pi \nu^*)] - \mathbb{E}_{\beta\pi} [\nu_{\mathcal{F}}^* - \nu^*] \\
&\leq \kappa \|\nu_{\mathcal{F}}^* - \nu^*\|_{\mathcal{D},1} + \kappa \|\mathcal{B}^\pi(\nu_{\mathcal{F}}^* - \nu^*)\|_{\mathcal{D},1} + \|\nu_{\mathcal{F}}^* - \nu^*\|_{\beta\pi,1} \\
&\leq \max \left(\kappa + \kappa \|\mathcal{B}^\pi\|_{\mathcal{D},1}, 1 \right) \left(\|\nu_{\mathcal{F}}^* - \nu^*\|_{\mathcal{D},1} + \|\nu_{\mathcal{F}}^* - \nu^*\|_{\beta\pi,1} \right) \\
&\leq \max \left(\kappa + \kappa \|\mathcal{B}^\pi\|_{\mathcal{D},1}, 1 \right) \cdot \epsilon_{approx}(\mathcal{F}),
\end{aligned}$$

530 where $\epsilon_{approx}(\mathcal{F}) := \sup_{\nu \in S \times A \rightarrow \mathbb{R}} \inf_{\nu_{\mathcal{F}} \in \mathcal{F}} \left(\|\nu_{\mathcal{F}} - \nu\|_{\mathcal{D},1} + \|\nu_{\mathcal{F}} - \nu\|_{\beta\pi,1} \right)$, due to the approxi-
531 mation with \mathcal{F} for ν .

532 For the term $L(\nu_{\mathcal{F}}^*) - J(\nu_{\mathcal{F}}^*)$, we have by definition that

$$L(\nu_{\mathcal{F}}^*) - J(\nu_{\mathcal{F}}^*) = \max_{\zeta \in \mathcal{H}} J(\nu_{\mathcal{F}}^*, \zeta) - \max_{\zeta \in S \times A \rightarrow \mathbb{R}} J(\nu_{\mathcal{F}}^*, \zeta) \leq 0$$

533 For the term $L(\hat{\nu}^*) - L(\nu_{\mathcal{F}}^*)$,

$$\begin{aligned}
L(\hat{\nu}^*) - L(\nu_{\mathcal{F}}^*) &= L(\hat{\nu}^*) - \hat{L}(\hat{\nu}^*) + \hat{L}(\hat{\nu}^*) - \hat{L}(\nu_{\mathcal{F}}^*) + \hat{L}(\nu_{\mathcal{F}}^*) - L(\nu_{\mathcal{F}}^*) \\
&\leq L(\hat{\nu}^*) - \hat{L}(\hat{\nu}^*) + \hat{L}(\nu_{\mathcal{F}}^*) - L(\nu_{\mathcal{F}}^*) \\
&\leq 2 \sup_{\nu \in \mathcal{F}} \left| L(\nu) - \hat{L}(\nu) \right| \\
&= 2 \sup_{\nu \in \mathcal{F}} \left| \max_{\zeta \in \mathcal{H}} J(\nu, \zeta) - \max_{\zeta \in \mathcal{H}} \hat{J}(\nu, \zeta) \right| \\
&\leq 2 \sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \hat{J}(\nu, \zeta) - J(\nu, \zeta) \right| \\
&= 2 \cdot \epsilon_{est}(\mathcal{F}),
\end{aligned}$$

534 where in the first inequality we have used the fact that $\hat{L}(\hat{\nu}^*) - \hat{L}(\nu_{\mathcal{F}}^*) \leq 0$ due to the optimality of

535 $\hat{\nu}^*$, and in the last step $\epsilon_{est}(\mathcal{F}) := \sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \hat{J}(\nu, \zeta) - J(\nu, \zeta) \right|$.

536 For the term $J(\hat{\nu}^*) - L(\hat{\nu}^*)$, we have

$$\begin{aligned} J(\hat{\nu}^*) - L(\hat{\nu}^*) &= \max_{\zeta \in S \times A \rightarrow \mathbb{R}} J(\hat{\nu}^*, \zeta) - \max_{\zeta \in \mathcal{H}} J(\hat{\nu}^*, \zeta) \\ &\leq \left(L + \frac{1+\gamma}{1-\gamma} C \right) \underbrace{\|\zeta_{\mathcal{H}}^* - \zeta^*\|_{\mathcal{D},1}}_{\leq \epsilon_{approx}(\mathcal{H})}, \end{aligned}$$

537 where $\epsilon_{approx}(\mathcal{H}) := \sup_{\zeta \in S \times A \rightarrow \mathbb{R}} \inf_{\zeta \in \mathcal{H}} \left(\|\zeta_{\mathcal{H}} - \zeta\|_{\mathcal{D},1} + \|\zeta_{\mathcal{H}} - \zeta\|_{\beta\pi,1} \right)$, due to the approxi-
538 mation with \mathcal{H} for ζ .

539 Finally, we can decompose the squared error as

$$\begin{aligned} &\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) (s, a) \cdot \hat{r}(s, a) \right] - \rho(\pi) \right)^2 \\ &\leq \frac{16C_r^2}{\eta} \left(\max \left(\kappa + \kappa \|\mathcal{B}^{\pi}\|_{\mathcal{D},1}, 1 \right) \epsilon_{approx}(\mathcal{F}) + \left(L + \frac{1+\gamma}{1-\gamma} C \right) \epsilon_{approx}(\mathcal{H}) \right) \\ &\quad + 4\epsilon_r + 8\epsilon_{stat} + \frac{32C_r^2}{\eta} \epsilon_{est}(\mathcal{F}) + 4\hat{\epsilon}_{opt}. \quad (24) \end{aligned}$$

Remark (Dual OPE estimator): We now reconcile the above derivations with the use of $\hat{\zeta}(s, a)$ as estimates of $w_{\pi/\mathcal{D}}(s, a)$. Note that in the implementation of DualDICE we use the estimator,

$$\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\hat{\zeta}(s, a) \cdot r \right]$$

540 for off-policy policy evaluation. In this case, the error can be decomposed as

$$\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\hat{\zeta}(s, a) \cdot r \right] - \mathbb{E}_{d^{\mathcal{D}}} \left[w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a) \right] \right)^2 \quad (25)$$

$$\leq 2 \left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\hat{\zeta}(s, a) \cdot r \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) (s, a) \cdot r \right] \right)^2 \quad (26)$$

$$+ 2 \left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) (s, a) \cdot r \right] - \mathbb{E}_{d^{\mathcal{D}}} \left[w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a) \right] \right)^2. \quad (27)$$

541 The second term above is the same as given in equation 19. The first term can be rewritten as,

$$\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\hat{\zeta}(s, a) \cdot r \right] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) (s, a) \cdot r \right] \right)^2 \leq C_r^2 \left\| \hat{\zeta} - \left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) \right\|_{\hat{\mathcal{D}}}^2,$$

542 which can be bounded as follows:

$$\begin{aligned} &\left\| \hat{\zeta} - \left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) \right\|_{\hat{\mathcal{D}}}^2 \\ &= \left\| \hat{\zeta} - \hat{\zeta}^* + \hat{\zeta}^* - \left(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^* \right) + \left(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^* \right) - \left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) \right\|_{\hat{\mathcal{D}}}^2 \\ &\leq 4 \left\| \hat{\zeta} - \hat{\zeta}^* \right\|_{\hat{\mathcal{D}}}^2 + 4 \left\| \left(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^* \right) - \left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) \right\|_{\hat{\mathcal{D}}}^2 + 4 \left\| \hat{\zeta}^* - \left(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^* \right) \right\|_{\hat{\mathcal{D}}}^2 \end{aligned} \quad (28)$$

543 where the first two terms correspond to optimization error $\hat{\epsilon}_{opt}$, and the last to approximation error
544 due to parametrization.

545 Specifically, when the output of our algorithm $\hat{\zeta}(s, a) = \left(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu} \right) (s, a)$ for $\forall (s, a) \in \hat{\mathcal{D}}$, the
546 extra term vanishes, and the error is the same as in equation 19.

547 D.2 Statistical Error

548 We analyze the statistical error ϵ_r , ϵ_{stat} and $\epsilon_{est}(\mathcal{F})$ in this section. We discussed in batch learning
549 setting with *i.i.d.* samples [48]. However, by exploiting blocking technique in Proposition 15 of [53],
550 following [2, 23, 9], all the sample complexity we provided can be easily generalized for single
551 β -mixing sample path, *i.e.*, $\{s_i, a_i, r_i, s'_i\}_{i=1}^N$ is strictly stationary and mixing in an exponential rate
552 with parameter $b, \chi > 0$ if $\beta_m = \mathcal{O}(\exp(-bm^{-\chi}))$, which we omit for the sake of exposition
553 simplicity.

554 **Bounding ϵ_r .** Recall that $\bar{R}(s, a) = \mathbb{E}_{\cdot|s,a}[r]$, so

$$\begin{aligned} \mathbb{E}[\epsilon_r] &\leq \left(\frac{1+\gamma}{1-\gamma}\right)^2 C^2 \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N r_i - \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N r_i \right] \right)^2 \right] \\ &= \left(\frac{1+\gamma}{1-\gamma}\right)^2 C^2 \mathbb{V} \left(\frac{1}{N} \sum_{i=1}^N r_i \right) \\ &\leq \frac{1}{N} \left(\frac{1+\gamma}{1-\gamma}\right)^2 C^2 \sup_{s,a} \mathbb{V}(r|s, a) = \mathcal{O} \left(\frac{1}{N} \right). \end{aligned} \quad (29)$$

555 Since $r(s, a)$ and $\hat{r}(s, a)$ is bounded, we can also obtain high-probability deviation bounds using
556 standard concentration inequalities [5].

Bounding $\epsilon_{est}(\mathcal{F})$. By definition, we have

$$\epsilon_{est}(\mathcal{F}) = \sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \hat{J}(\nu, \zeta) - J(\nu, \zeta) \right|,$$

557 which can be bounded using a covering-number argument outlined below.

558 We will need Pollard's tail inequality that relates maximum deviation to the covering number of a
559 function class:

560 **Lemma 4.** [35] Let \mathcal{G} be a permissible class of $\mathcal{Z} \rightarrow [-M, M]$ functions and $\{Z_i\}_{i=1}^N$ are i.i.d. sam-
561 ples from some distribution. Then, for any given $\epsilon > 0$,

$$\mathbb{P} \left(\sup_{g \in \mathcal{G}} \left| \frac{1}{N} \sum_{i=1}^N g(Z_i) - \mathbb{E}[g(Z)] \right| > \epsilon \right) \leq 8 \mathbb{E} \left[\mathcal{N}_1 \left(\frac{\epsilon}{8}, \mathcal{G}, \{Z_i\}_{i=1}^N \right) \right] \exp \left(\frac{-N\epsilon^2}{512M^2} \right).$$

562 The covering number can then be bounded in terms of the function class's pseudo-dimension:

Lemma 5. [Corollary 3, [20]] For any set \mathcal{X} , any points $x^{1:N} \in \mathcal{X}^N$, any class \mathcal{F} of functions on \mathcal{X} taking values in $[0, M]$ with pseudo-dimension $D_{\mathcal{F}} < \infty$, and any $\epsilon > 0$,

$$\mathcal{N}_1(\epsilon, \mathcal{F}, x^{1:N}) \leq e(D_{\mathcal{F}} + 1) \left(\frac{2eM}{\epsilon} \right)^{D_{\mathcal{F}}}.$$

563 With the above technical lemmas, we are ready to bound $\epsilon_{est}(\mathcal{F})$.

Lemma 6 (Statistical error $\epsilon_{est}(\mathcal{F})$). Under Assumption 1, if f^* is L -Lipschitz continuous, with at least probability $1 - \delta$,

$$\epsilon_{est}(\mathcal{F}) = \mathcal{O} \left(\sqrt{\frac{\log N + \log \frac{1}{\delta}}{N}} \right).$$

564 *Proof.* Denote $h_{\nu, \zeta}(s, a, s', a', s_0, a_0) = (\nu(s, a) - \gamma\nu(s', a'))\zeta(s, a) - f^*(\zeta(s, a)) - (1 -$
565 $\gamma)\nu(s_0, a_0)$, we use lemma 4 with $\mathcal{Z} = \underbrace{S \times A \times S \times A}_{d^D \pi} \times \underbrace{S \times A}_{\beta \pi}$, $Z_i = (s_i, a_i, s'_i, a'_i, s_0^i, a_0^i)$

566 and $\mathcal{G} = h_{\mathcal{F} \times \mathcal{H}}$.

567 We first show that $\forall h_{\nu, \zeta} \in \mathcal{G}$ is bounded. Recall $\nu \in \mathcal{F}$ and $\zeta \in \mathcal{H}$ are bounded by $\frac{1}{1-\gamma}C$ and C ,
568 then, $h_{\nu, \zeta}$ will be bounded by $M_1 = \frac{1+\gamma}{1-\gamma}C^2 + (1+L)C + |f^*(0)|$. Specifically,

$$\begin{aligned} \|h_{\nu, \zeta}\|_{\infty} &\leq (1+\gamma)\|\nu\|_{\infty}\|\zeta\|_{\infty} + (1-\gamma)\|\nu\|_{\infty} + \|f^*(\zeta)\|_{\infty} \\ &\leq \frac{1+\gamma}{1-\gamma}C^2 + C + \|f^*(\zeta) - f^*(0)\|_{\infty} + |f^*(0)| \\ &\leq \frac{1+\gamma}{1-\gamma}C^2 + C + L\|\zeta\|_{\infty} + |f^*(0)| \\ &\leq \frac{1+\gamma}{1-\gamma}C^2 + C + LC + |f^*(0)|. \end{aligned}$$

569 Thus,

$$\begin{aligned} \mathbb{P} \left(\sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \hat{J}(\nu, \zeta) - J(\nu, \zeta) \right| \geq \epsilon \right) &= \mathbb{P} \left(\sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \frac{1}{N} \sum_{i=1}^N h_{\nu, \zeta}(Z_i) - \mathbb{E}[h_{\nu, \zeta}] \right| \geq \epsilon \right) \\ &\leq 8\mathbb{E} \left[\mathcal{N}_1 \left(\frac{\epsilon}{8}, \mathcal{G}, \{Z_i\}_{i=1}^N \right) \right] \exp \left(\frac{-N\epsilon^2}{512M_1^2} \right). \end{aligned} \quad (30)$$

570 We bound the distance in \mathcal{G} ,

$$\begin{aligned} &\frac{1}{N} \sum_{i=1}^N |h_{\nu_1, \zeta_1}(Z_i) - h_{\nu_2, \zeta_2}(Z_i)| \\ &\leq \frac{\left(L + \frac{1+\gamma}{1-\gamma} C \right)}{N} \sum_{i=1}^N |\zeta_1(s_i, a_i) - \zeta_2(s_i, a_i)| + \frac{C}{N} \sum_{i=1}^N |\nu_1(s_i, a_i) - \nu_2(s_i, a_i)| \\ &\quad + \frac{\gamma C}{N} \sum_{i=1}^N |\nu_1(s'_i, a'_i) - \nu_2(s'_i, a'_i)| + \frac{(1-\gamma)}{N} \sum_{i=1}^N |\nu_1(s_0^i, a_0^i) - \nu_2(s_0^i, a_0^i)|, \end{aligned}$$

571 which leads to

$$\begin{aligned} &\mathcal{N}_1 \left(\left(L + \frac{2+\gamma-\gamma^2}{1-\gamma} C + (1-\gamma) \right) \epsilon', \mathcal{G}, \{Z_i\}_{i=1}^N \right) \\ &\leq \mathcal{N}_1 \left(\epsilon', \mathcal{H}, \{s_i, a_i\}_{i=1}^N \right) \mathcal{N}_1 \left(\epsilon', \mathcal{F}, \{s_i, a_i\}_{i=1}^N \right) \mathcal{N}_1 \left(\epsilon', \mathcal{F}, \{s'_i, a'_i\}_{i=1}^N \right) \mathcal{N}_1 \left(\epsilon', \mathcal{F}, \{s_0^i, a_0^i\}_{i=1}^N \right). \end{aligned} \quad (31)$$

Applying lemma 5, we can bound the covering number. Denote the pseudo-dimension of \mathcal{F} and \mathcal{H} as D_ν and D_ζ , then, we have

$$\mathcal{N}_1 \left(\left(L + \frac{2+\gamma-\gamma^2}{1-\gamma} C + (1-\gamma) \right) \epsilon', \mathcal{G}, \{Z_i\}_{i=1}^N \right) \leq e^4 (D_\mathcal{F} + 1)^3 (D_\mathcal{H} + 1) \left(\frac{4eM_1}{\epsilon'} \right)^{3D_\mathcal{F} + D_\mathcal{H}},$$

572 which implies

$$\begin{aligned} &\mathcal{N}_1 \left(\frac{\epsilon}{8}, \mathcal{G}, \{Z_i\}_{i=1}^N \right) \\ &\leq e^4 (D_\mathcal{F} + 1)^3 (D_\mathcal{H} + 1) \left(\frac{32 \left(L + \frac{2+\gamma-\gamma^2}{1-\gamma} C + (1-\gamma) \right) eM_1}{\epsilon} \right)^{3D_\mathcal{F} + D_\mathcal{H}} := C_1 \left(\frac{1}{\epsilon} \right)^{D_1}, \end{aligned} \quad (32)$$

573 where $C_1 = e^4 (D_\mathcal{F} + 1)^3 (D_\mathcal{H} + 1) \left(32 \left(L + \frac{2+\gamma-\gamma^2}{1-\gamma} C + (1-\gamma) \right) eM_1 \right)^{D_1}$ and $D_1 = 3D_\mathcal{F} +$
574 $D_\mathcal{H}$.

Combine this result with equation 30, we immediately obtain the statistical error, *i.e.*,

$$\mathbb{P} \left(\sup_{\nu \in \mathcal{F}, \zeta \in \mathcal{H}} \left| \hat{J}(\nu, \zeta) - J(\nu, \zeta) \right| \geq \epsilon \right) \leq 8C_1 \left(\frac{1}{\epsilon} \right)^{D_1} \exp \left(\frac{-N\epsilon^2}{512M_1^2} \right).$$

By setting $\epsilon = \sqrt{\frac{C_2(\log N + \log \frac{1}{\delta})}{N}}$ with $C_2 = \max \left((8C_1)^{\frac{2}{D_1}}, 512M_1D_1, 512M_1, 1 \right)$, we have

$$8C_1 \left(\frac{1}{\epsilon} \right)^{D_1} \exp \left(\frac{-N\epsilon^2}{512M_1^2} \right) \leq \delta.$$

575

□

576 **Bounding ϵ_{stat} .** As $\hat{\nu}^*$ is a random variable, we need to bound the following instead:

$$\begin{aligned} & \sqrt{\epsilon_{stat}} \\ &= \left| \hat{\mathbb{E}}_{s,a,s',a'} [(\hat{\nu}^*(s,a) - \gamma \hat{\nu}^*(s',a')) r(s,a)] - \mathbb{E}_{s,a,s',a'} [(\hat{\nu}^*(s,a) - \gamma \hat{\nu}^*(s',a')) r(s,a)] \right| \\ &\leq \sup_{\nu \in \mathcal{F}} \left| \hat{\mathbb{E}}_{s,a,s',a'} [(\nu(s,a) - \gamma \nu(s',a')) r(s,a)] - \mathbb{E}_{s,a,s',a'} [(\nu(s,a) - \gamma \nu(s',a')) r(s,a)] \right|, \end{aligned}$$

577 which can be done using a similar argument as above.

Lemma 7 (Statistical error ϵ_{stat}). *Under Assumption 1, with at least probability $1 - \delta$,*

$$\epsilon_{stat} = \mathcal{O} \left(\frac{\log N + \log \frac{1}{\delta}}{N} \right).$$

Proof. We first show that $\forall \nu \in \mathcal{H}$, $(\nu(s,a) - \gamma \nu(s',a')) r(s,a)$ is bounded by $M_2 = \frac{1+\gamma}{1-\gamma} C^2$, i.e.,

$$\|(\nu - \gamma \nu') \cdot r\|_\infty \leq (1 + \gamma) C \|\nu\|_\infty \leq \frac{1 + \gamma}{1 - \gamma} C^2.$$

578 Then, we apply the lemma 4 with $\mathcal{Z} = S \times A \times S \times A$, $Z_i = (s_i, a_i, s'_i, a'_i)$, and $\mathcal{G} = (\nu - \gamma \nu) \cdot r$,

$$\mathbb{P} \left(\sup_{\nu \in \mathcal{F}} \left| \hat{\mathbb{E}}_{\mathcal{Z}} [(\nu - \hat{\mathcal{B}}^\pi \nu) \cdot r] - \mathbb{E} [(\nu - \mathcal{B}^\pi \nu) \cdot r] \right| \geq \epsilon \right) \quad (33)$$

$$\leq 8 \mathbb{E} \left[\mathcal{N}_1 \left(\frac{\epsilon}{8}, \mathcal{G}, \{Z_i\}_{i=1}^N \right) \right] \exp \left(\frac{-N \epsilon^2}{512 M_2^2} \right). \quad (34)$$

579 Similarly, we have

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N |(\nu_1 - \gamma \nu_1) \cdot r(Z_i) - (\nu_2 - \gamma \nu_2) \cdot r(Z_i)| \\ &\leq \frac{C}{N} \sum_{i=1}^N |\nu_1(s_i, a_i) - \nu_2(s_i, a_i)| + \frac{\gamma C}{N} |\nu_1(s'_i, a'_i) - \nu_2(s'_i, a'_i)|, \end{aligned}$$

580 leading to

$$\mathcal{N}_1 \left((1 + \gamma) C \epsilon', \mathcal{G}, \{Z_i\}_{i=1}^N \right) \leq \mathcal{N}_1 \left(\epsilon', \mathcal{F}, \{s_i, a_i\}_{i=1}^N \right) \mathcal{N}_1 \left(\epsilon', \mathcal{F}, \{s'_i, a'_i\}_{i=1}^N \right). \quad (35)$$

581 Applying lemma 5, we bound the covering number as

$$\mathcal{N}_1 \left((1 + \gamma) C \epsilon', \mathcal{G}, \{Z_i\}_{i=1}^N \right) \leq e^2 (D_{\mathcal{F}} + 1)^2 \left(\frac{2e M_2}{\epsilon'} \right)^{2D_{\mathcal{F}}}, \quad (36)$$

which implies

$$\mathcal{N}_1 \left(\frac{\epsilon}{8}, \mathcal{G}, \{Z_i\}_{i=1}^N \right) \leq e^2 (D_{\mathcal{F}} + 1)^2 \left(\frac{16(1 + \gamma) C e M_2}{\epsilon} \right)^{2D_{\mathcal{F}}} := C_3 \left(\frac{1}{\epsilon} \right)^{D_2},$$

582 with $C_3 := e^2 (D_{\mathcal{F}} + 1)^2 (16(1 + \gamma) C e M_2)^{D_2}$ and $D_2 = 2D_{\mathcal{F}}$.

583 We achieve the statistical error bound, i.e.,

$$\mathbb{P}(\sqrt{\epsilon_{stat}} \geq \epsilon) \leq 8 C_3 \left(\frac{1}{\epsilon} \right)^{D_2} \exp \left(\frac{-N \epsilon^2}{512 M_2^2} \right). \quad (37)$$

By setting $\epsilon = \sqrt{\frac{C_4 (\log N + \log \frac{1}{\delta})}{N}}$ with $C_4 = \max \left((8 C_3)^{\frac{2}{D_2}}, 512 M_2 D_2, 512 M_2, 1 \right)$, we have

$$8 C_3 \left(\frac{1}{\epsilon} \right)^{D_2} \exp \left(\frac{-N \epsilon^2}{512 M_2^2} \right) \leq \delta.$$

584 Therefore, we have $\epsilon_{stat} = \mathcal{O} \left(\frac{\log N + \log \frac{1}{\delta}}{N} \right)$, with $1 - \delta$ probability. \square

585 D.3 Putting It All Together

586 **Theorem 2** Under Assumptions 1 and 3, with $f(x) = \frac{1}{2}x^2$, the mean squared error of DualDICE's
 587 estimate is bounded by

$$\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [\hat{\zeta}(s, a) \cdot r] - \rho(\pi) \right)^2 \right] = \tilde{\mathcal{O}} \left(\epsilon_{approx}(\mathcal{F}, \mathcal{H}) + \epsilon_{opt} + \frac{1}{\sqrt{N}} \right),$$

588 where $\mathbb{E}[\cdot]$ is taken w.r.t. randomness both in the sampling of $\mathcal{D} \sim d^{\mathcal{D}}$ and in the algorithm, $\tilde{\mathcal{O}}(\cdot)$
 589 ignores logarithmic factors, and the error terms are defined in equation 40.

590 *Proof.* By equations 25 and 28, the error can be decomposed as

$$\begin{aligned} & \mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [\hat{\zeta}(s, a) \cdot r] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)] \right)^2 \right] \\ & \leq 2\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [\hat{\zeta}(s, a) \cdot r] - \hat{\mathbb{E}}_{d^{\mathcal{D}}} [(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot r] \right)^2 \right] \\ & \quad + 2\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot r] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot \bar{R}(s, a)] \right)^2 \right] \\ & \leq 8C_r^2 \mathbb{E} \left[\left(\|\hat{\zeta} - \hat{\zeta}^*\|_{\hat{\mathcal{D}}}^2 + \|(\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*) - (\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})\|_{\hat{\mathcal{D}}}^2 \right) \right] + 8C_r^2 \mathbb{E} \left[\|\hat{\zeta}^* - (\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)\|_{\hat{\mathcal{D}}}^2 \right] \\ & \quad + 2\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \hat{r}(s, a)] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot r(s, a)] \right)^2 \right]. \end{aligned} \quad (38)$$

We can bound the last term, $\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \hat{r}(s, a)] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot r(s, a)] \right)^2 \right]$,
 by straightforwardly combining equation 29, lemma 6 and lemma 7 into equation 24. Specifically, by
 lemma 6, we have

$$\mathbb{E} [\epsilon_{est}(\mathcal{F})] = \sqrt{\frac{C_2 \log N + \log \frac{1}{\delta_1}}{N}} (1 - \delta_1) + 2\delta_1 M_1 = \mathcal{O} \left(\sqrt{\frac{\log N}{N}} \right),$$

by setting $\delta_1 = \frac{1}{\sqrt{N}}$. Similarly, we have

$$\mathbb{E} [\epsilon_{stat}] = \frac{C_4 \left(\log N + \log \frac{1}{\delta_2} \right)}{N} (1 - \delta_2) + 2\delta_2 M_2 = \mathcal{O} \left(\frac{\log N}{N} \right),$$

591 where the last equation comes from by setting $\delta_2 = \frac{1}{N}$. Plug these results into equation 24, we have

$$\begin{aligned} & \mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} [(\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu})(s, a) \cdot \hat{r}(s, a)] - \mathbb{E}_{d^{\mathcal{D}}} [w_{\pi/\mathcal{D}}(s, a) \cdot r(s, a)] \right)^2 \right] \\ & \leq \mathcal{O}(\epsilon_{approx}(\mathcal{F}) + \epsilon_{approx}(\mathcal{H}) + \epsilon_{opt}) + \tilde{\mathcal{O}} \left(\frac{1}{N} + \sqrt{\frac{1}{N}} \right), \end{aligned} \quad (39)$$

592 where $\epsilon_{opt} = \mathbb{E}[\hat{\epsilon}_{opt}]$, $\epsilon_{approx}(\mathcal{F}) := \sup_{\nu \in S \times A \rightarrow \mathbb{R}} \inf_{\nu \in \mathcal{F}} \left(\|\nu_{\mathcal{F}} - \nu\|_{\mathcal{D}, 1} + \|\nu_{\mathcal{F}} - \nu\|_{\beta\pi, 1} \right)$, and
 593 $\epsilon_{approx}(\mathcal{H}) := \sup_{\zeta \in S \times A \rightarrow \mathbb{R}} \inf_{\zeta \in \mathcal{H}} \left(\|\zeta_{\mathcal{H}} - \zeta\|_{\mathcal{D}, 1} + \|\zeta_{\mathcal{H}} - \zeta\|_{\beta\pi, 1} \right)$, due to the approximation
 594 with \mathcal{F} for ν and \mathcal{H} for ζ , respectively.

595 The first term in equation 38, $\mathbb{E} \left[\left(\|\hat{\zeta} - \hat{\zeta}^*\|_{\hat{\mathcal{D}}}^2 + 8 \left\| (\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*) - (\hat{\nu} - \hat{\mathcal{B}}^{\pi} \hat{\nu}) \right\|_{\hat{\mathcal{D}}}^2 \right) \right]$, is also the
 596 optimization error ϵ_{opt} .

597 The second term, $\mathbb{E} \left[\|\hat{\zeta}^* - (\hat{\nu}^* - \hat{\mathcal{B}}^{\pi} \hat{\nu}^*)\|_{\hat{\mathcal{D}}}^2 \right]$, is due to the parametrization by \mathcal{F} and \mathcal{H} .

598 Define the approximation error

$$\epsilon_{approx}(\mathcal{F}, \mathcal{H}) = \epsilon_{approx}(\mathcal{F}) + \epsilon_{approx}(\mathcal{H}) + \mathbb{E} \left[\left\| \hat{\zeta}^* - \left(\hat{\nu}^* - \hat{\mathcal{B}}^\pi \hat{\nu}^* \right) \right\|_{\hat{\mathcal{D}}}^2 \right], \quad (40)$$

combine equation 39 and the extra errors, we immediately have

$$\mathbb{E} \left[\left(\hat{\mathbb{E}}_{d^{\mathcal{D}}} \left[\hat{\zeta}(s, a) \cdot \hat{r}(s, a) \right] - \mathbb{E}_{d^{\mathcal{D}}} \left[w_{\pi/\mathcal{D}}(s, a) \cdot r(s, a) \right] \right)^2 \right] = \tilde{\mathcal{O}} \left(\epsilon_{approx}(\mathcal{F}, \mathcal{H}) + \epsilon_{opt} + \frac{1}{\sqrt{N}} \right),$$

599 which is the first conclusion.

600

□

601 D.4 Optimization Error

602 In this section, we characterize the optimization error $\hat{\epsilon}_{opt}$. With different parametrizations for $(\mathcal{F}, \mathcal{H})$
 603 and different optimization algorithms for $\hat{J}(\nu, \zeta)$, the convergence rate of ϵ_1 will be different. For
 604 general parametrization of $(\mathcal{F}, \mathcal{H})$ as neural network, how to quantitatively analyze the optimization
 605 error is still an open problem and out of the scope of this paper. We focus on the tabular, linear or
 606 kernel parametrization for $(\mathcal{F}, \mathcal{H})$. Let $(\mathcal{F}, \mathcal{H})$ are the family of linear models with basis function
 607 $\psi(s, a) \in \mathbb{R}^p$. The tabular and kernel version can be easily generalized by treating ψ as indicator
 608 vectors or infinite dimension feature mapping, respectively, and we omit here. Then, we can
 609 parametrize $\nu(s, a) = w_\nu^\top \psi(s, a)$ and $\zeta(s, a) = w_\zeta^\top \psi(s, a)$ with $w_\nu, w_\zeta \in \mathbb{R}^p$. Then, the
 610 optimization reduces to

$$\min_{w_\nu \in \mathcal{F}} \max_{w_\zeta \in \mathcal{H}} \hat{J}(w_\nu, w_\zeta) := w_\nu^\top \mathcal{A} w_\zeta - \frac{1}{N} \sum_{i=1}^N f^*(w_\zeta^\top \psi(s_i, a_i)_{\mathcal{H}}) - w_\nu^\top b, \quad (41)$$

611 where $\mathcal{A} = \frac{1}{N} \sum_{i=1}^N (\psi(s_i, a_i) - \gamma \psi(s'_i, a'_i)) \psi^\top(s_i, a_i) \in \mathbb{R}^{p \times p}$ and $b = \frac{(1-\gamma)}{N} \sum_{i=1}^N \psi(s_0^i, a_0^i)$.

612 We have

$$\begin{aligned} \hat{\epsilon}_{opt} &= \left\| \hat{\zeta} - \hat{\zeta}^* \right\|_{\hat{\mathcal{D}}}^2 + \left\| \left(\hat{\nu}^* - \hat{\mathcal{B}}^\pi \hat{\nu}^* \right) - \left(\hat{\nu} - \hat{\mathcal{B}}^\pi \hat{\nu} \right) \right\|_{\hat{\mathcal{D}}}^2 \\ &\leq \left\| \Psi \right\|_2^2 \left\| \hat{w}_\zeta - \hat{w}_\zeta^* \right\|^2 + \left\| \Phi \right\|_2^2 \left\| \hat{w}_\nu - \hat{w}_\nu^* \right\|^2 \\ &\leq \max \left(\left\| \Psi \right\|_2^2 + \left\| \Phi \right\|_2^2 \right) \left(\left\| \hat{w}_\zeta - \hat{w}_\zeta^* \right\|^2 + \left\| \hat{w}_\nu - \hat{w}_\nu^* \right\|^2 \right), \end{aligned} \quad (42)$$

613 where $\Psi = [\psi(s_i, a_i)]_{i=1}^N \in \mathbb{R}^{N \times p}$ and $\Phi = [\psi(s_i, a_i) - \gamma \psi(s'_i, a'_i)]_{i=1}^N \in \mathbb{R}^{N \times p}$.

614 In general case, the optimization 41 is convex-concave, therefore, the vanilla stochastic gradient
 615 descent converges in rate $\mathcal{O} \left(\frac{1}{\sqrt{T}} \right)$ in terms of the primal-dual gap. Specifically, we have $f(x) = \frac{1}{2} x^2$,
 616 which will lead $\frac{1}{N} \sum_{i=1}^N f^*(w_\zeta^\top \psi(s_i, a_i)_{\mathcal{H}}) = \|w_\zeta\|_{\mathcal{C}}^2$ with $\mathcal{C} = \frac{1}{N} \sum_{i=1}^N \psi(s_i, a_i) \psi^\top(s_i, a_i) \in$
 617 $\mathbb{R}^{d \times d}$. Under the assumption as [11],

618 **Assumption 8.** \mathcal{A} has full rank, \mathcal{C} is strictly positive definite, and the feature vector $\psi(s, a)$ is
 619 uniformly bounded.

620 We discuss the optimization error $\epsilon_{opt} := \mathbb{E}[\epsilon_1]$, where the $\mathbb{E}[\cdot]$ w.r.t. the randomness in the algorithm,
 621 in two algorithms for equation 41,

622 • **SVRG** We can easily verify that the T -step solution of SVRG, $(\hat{\nu}_T, \hat{\zeta}_T)$, converges to
 623 $(\hat{\nu}^*, \hat{\zeta}^*)$ in linear rate $\mathcal{O}(\exp(-T))$ in terms of $\mathbb{E} \left[\left\| \hat{w}_\nu^T - \hat{w}_\nu^* \right\|^2 + \left\| \hat{w}_\zeta^T - \hat{w}_\zeta^* \right\|^2 \right]$ follow-
 624 ing [11], where the expectation w.r.t. the randomness in the SVRG. Specifically, we have
 625

$$\hat{\epsilon}_{opt} = \mathcal{O}(\exp(-T)). \quad (43)$$

626 • **SGD** Although the optimization equation 41 is not strongly convex-concave, we can still
 627 prove $\mathcal{O} \left(\frac{1}{T} \right)$ convergence rate.

628
629

Lemma 9. *Let the stepsize τ_t decay in $\mathcal{O}(\frac{1}{t})$, assume the norm of the stochastic gradient is bounded, under Assumption 8, we have*

$$\hat{\epsilon}_{opt} = \mathcal{O}\left(\frac{1}{T}\right). \quad (44)$$

Proof. Denote $\hat{\theta}_t = \left[\hat{w}_\nu^t, \frac{1}{\sqrt{\rho}}\hat{w}_\theta^t\right]$ and

$$G_t = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{\rho}} \end{bmatrix} \cdot \underbrace{\begin{bmatrix} 0 & \sqrt{\rho}\hat{\mathcal{A}}_t \\ -\sqrt{\rho}\hat{\mathcal{A}}_t^\top & \rho\hat{\mathcal{C}}_t \end{bmatrix}}_{\hat{Q}} \cdot \begin{bmatrix} w_\nu^t \\ \frac{1}{\sqrt{\rho}}w_\zeta^t \end{bmatrix} - \begin{bmatrix} \hat{b}_t \\ 0 \end{bmatrix}$$

630

as the unbiased stochastic gradient with $\mathbb{E}[G_t] = g_t$, we have the update rule

631

as $\theta_{t+1} = \theta_t - \Sigma_t G_t$, $\Sigma_t = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{\rho}} \end{bmatrix} \sigma_t$ We denote $\delta_t = \frac{1}{2} \left\| \hat{\theta}_t - \hat{\theta}^* \right\|^2 =$

632

$\frac{1}{2} \left[\left\| \hat{w}_\nu^t - \hat{w}_\nu^* \right\|^2 + \frac{1}{\rho} \left\| \hat{w}_\zeta^t - \hat{w}_\zeta^* \right\|^2 \right]$ and $\Delta_t = \mathbb{E}[\delta_t]$. Then, we have

$$\delta_{t+1} = \frac{1}{2} \left\| \hat{\theta}_t - \Sigma_t G_t - \hat{\theta}^* \right\|^2 \leq \delta_t + \frac{1}{2} \sigma_t^2 \left(1 + \frac{1}{\rho} \right) \|G_t\|^2 - \left(\hat{\theta}_t - \hat{\theta}^* \right)^\top (\Sigma_t G_t).$$

633

Take the expectation on both sides and $\mathbb{E}[\|G_t\|^2] \leq K^2$,

$$\Delta_{t+1} = \Delta_t + \frac{\sigma_t^2}{2} \left(1 + \frac{1}{\rho} \right) K^2 - \mathbb{E} \left[\left(\hat{\theta}_t - \hat{\theta}^* \right)^\top (\Sigma_t g_t) \right]. \quad (45)$$

634

As shown in [11], under Assumption 8 and set $\rho = \frac{8\lambda_{\max}(\mathcal{A}\mathcal{C}^{-1}\mathcal{A})}{\lambda_{\min}(\mathcal{C})}$, the matrix $Q := \mathbb{E}[\hat{Q}]$ has positive real eigenvalue and

635

$$\lambda_{\max}(Q) \leq 9 \frac{\lambda_{\max}(\mathcal{C})}{\lambda_{\min}(\mathcal{C})} \lambda_{\max}(\mathcal{A}\mathcal{C}^{-1}\mathcal{A}), \quad \lambda_{\min}(Q) \geq \frac{8}{9} \lambda_{\min}(\mathcal{A}\mathcal{C}^{-1}\mathcal{A}).$$

636

On the other hand, with the first-order optimality condition, we can show that $Q\hat{\theta}^* = \hat{b}$. Then, we have

637

$$\begin{aligned} \mathbb{E} \left[\left(\hat{\theta}_t - \hat{\theta}^* \right)^\top (\Sigma_t g_t) \right] &= \mathbb{E} \left[\left(\hat{\theta}_t - \hat{\theta}^* \right)^\top \Sigma_t^2 (Q\hat{\theta}_t - \hat{b}) \right] \\ &= \mathbb{E} \left[\left(\hat{\theta}_t - \hat{\theta}^* \right)^\top \Sigma_t^2 Q (\hat{\theta}_t - \hat{\theta}^*) \right] \geq 2\lambda_{\min}(Q) \left(1 + \frac{1}{\rho} \right) \sigma_t^2 \Delta_t. \end{aligned}$$

638

Plug this into the equation 30, we obtain the recursion,

$$\Delta_{t+1} \leq \Delta_t + \frac{\sigma_t^2}{2} \left(1 + \frac{1}{\rho} \right) K^2 - 2\lambda_{\min}(Q) \left(1 + \frac{1}{\rho} \right) \sigma_t^2 \Delta_t \leq (1 - 2c\sigma_t) \Delta_t + \frac{1 + \frac{1}{\rho}}{2} \sigma_t^2 K^2, \quad (46)$$

639

with $c = \lambda_{\min}(Q) \left(1 + \frac{1}{\rho} \right)$. By setting $\sigma_t > \frac{1}{2ct}$, $\Delta_T = \mathcal{O}(\frac{1}{T})$.

640

□

641

Using the above results for linear parametrization, we can reach the following corollary of Theorem 2.

642

Corollary 10. *Under the conditions of Theorem 2 and with linear parametrization of (ν, ζ) and under Assumption 8, after T -iteration, we have $\hat{\epsilon}_{opt} = \mathcal{O}(\exp(-T))$ for SVRG and $\hat{\epsilon}_{opt} = \mathcal{O}(\frac{1}{T})$ for SGD.*

644