# Beyond the Return: Off-policy Function Estimation under User-specified Error-measuring Distributions

Audrey Huang<sup>1</sup> Nan Jiang<sup>1</sup>

# Abstract

Off-policy evaluation often refers to two related tasks: estimating the expected return of a policy and estimating its value function (or other functions of interest, such as density ratios). While recent works on marginalized importance sampling (MIS) show that the former can enjoy provable guarantees under realizable function approximation, the latter is only known to be feasible under much stronger assumptions such as prohibitively expressive discriminators. In this work, we provide guarantees for off-policy function estimation under only realizability, by imposing proper regularization on the MIS objectives. Compared to commonly used regularization in MIS, our regularizer is much more flexible and can account for an arbitrary user-specified distribution, under which the learned function will be close to the groundtruth. We provide *exact* characterization of the optimal dual solution that needs to be realized by the discriminator class, which determines the data-coverage assumption in the case of value-function learning. As another surprising observation, the regularizer can be altered to relax the data-coverage requirement, and completely eliminate it in the ideal case with strong side information.

#### 1. Introduction

Off-policy evaluation (OPE) often refers to two related tasks in reinforcement learning (RL): estimating the expected return of a target policy using a dataset collected from a different *behavior* policy, versus estimating the policy's value function (or other functions of interest, such as density ratios). The former is crucial to hyperparameter tuning and verifying the performance of a policy before real-world deployment in offline RL (Voloshin et al., 2019; Paine et al., 2020; Zhang & Jiang, 2021). The latter, on the other hand, plays an important role in (both online and offline) training, often as the subroutine of actor-critic-style algorithms (Lagoudakis & Parr, 2003; Liu et al., 2019), but is also generally more difficult than the former: if an accurate value function is available, one could easily estimate the return by plugging in the initial distribution.

Between the two tasks, the theoretical nature of off-policy return estimation is relatively well understood, especially in terms of the function-approximation assumptions needed for sample-complexity guarantees. Among the available algorithms, importance sampling (IS) and its variants (Precup et al., 2000; Thomas et al., 2015; Jiang & Li, 2016) do not require any function approximation, but incur exponentialin-horizon variance. Fitted-Q Evaluation (Ernst et al., 2005; Le et al., 2019) can enjoy polynomial sample complexity under appropriate coverage assumptions, but the guarantee requires the function class to satisfy the strong Bellmancompleteness assumption, i.e. closure under the Bellman operator (Chen & Jiang, 2019; Xie et al., 2021). Marginalized importance sampling (MIS) methods, which have gained significant attention recently (Liu et al., 2018; Xie et al., 2019; Uehara et al., 2020; Nachum et al., 2019a), use two function classes to simultaneously approximate the value and the density-ratio (or weight) function and optimize minimax objectives. Notably, it is the only family of methods known to produce accurate return estimates with a polynomial sample complexity, when the function classes only satisfy the relatively weak realizability assumptions (i.e., they contain the true value and weight functions).

In comparison, little is known about off-policy function estimation, and the guarantees are generally less desirable. Not only do the limitations of IS and FQE on return estimation carry over to this more challenging task, but MIS also loses its major advantage over FQE: despite the somewhat misleading impression left by many prior works, that MIS can handle function estimation the same way as return estimation, <sup>1</sup> MIS for function estimation often requires unrealistic assumptions, such as prohibitively expressive discriminators. For concreteness, a typical guarantee for

<sup>&</sup>lt;sup>1</sup>Department of Computer Science, University of Illinois, Urbana-Champaign, USA. Correspondence to: <audreyh5@illinois.edu>, <nanjiang@illinois.edu>.

Decision Awareness in Reinforcement Learning Workshop at the 39th International Conference on Machine Learning (ICML), Baltimore, Maryland, USA, 2022. Copyright 2022 by the author(s)

<sup>&</sup>lt;sup>1</sup> For example, Liu et al. (2019) assume access to a weight estimation oracle and cited Liu et al. (2018) as a possible instance.

function estimation from MIS looks like the following (see e.g., Theorem 4 of Liu et al. (2018) and Lemmas 1 and 3 of Uehara et al. (2020)):

**Proposition 1** (Function estimation guarantee for MIS, informal). Suppose the offline data distribution  $d^D$  satisfies  $d^D(s, a) > 0, \forall s, a$ . Given value-function class Q with  $q^{\pi} \in Q$  and weight class  $\mathcal{W} = \mathbb{R}^{S \times \mathcal{A}}, q^{\pi} =$  $\arg \min_{q \in Q} \max_{w \in \mathcal{W}} L(w, q)$  for some appropriate population loss function L.

To enable the identification of the value function  $q^{\pi}$ , the result requires the discriminator class W to be the *space* of all possible functions over the state-action space ( $W = \mathbb{R}^{S \times A}$ ). In the finite-sample regime, using such a class incurs a sample complexity that depends on the size of the state-action space, which completely beats the purpose of function approximation.

In addition, these results only hold asymptotically, where the function of interest can be exactly identified in a pointwise manner. Such an overly strong guarantee is unrealistic in the finite-sample regime, where one can only hope to approximate the function well in an average sense under some distribution, i.e., finite-sample performance guarantees should ideally bound  $\|\widehat{q} - q^{\pi}\|_{2,\nu}$  for the learned  $\widehat{q}$ , where  $\|\cdot\|_{2,\nu}$  is  $\nu$ -weighted 2-norm. Such fine-grained analyses are non-existent in MIS. Even in the broader literature, such results not only require Bellman-completeness-type assumptions (Uehara et al., 2021), they also come with some fixed  $\nu$  (which is not necessarily  $d^D$ ; see Section 2) and the user has no freedom in choosing  $\nu$ . This creates a gap in the literature, as downstream learning algorithms that use offpolicy function estimation as a subroutine often assume the estimation to be accurate under certain specific distributions. For example, in the setting of online policy optimization, Abbasi-Yadkori et al. (2019) require value estimates to be accurate on the occupancy of the (unknown) optimal policy, and (Kakade & Langford, 2002) require them to be accurate on the occupancy of the learning policy at each iteration <sup>2</sup>. Choosing  $\nu$  to be the initial state-action distribution is well-suited for off-policy return estimation(see Appendix F). In offline RL, Liu et al. (2019) assume access to weight and value function estimation oracles on  $\nu = d^D$ .

To summarize, below are two important open problems on off-policy function estimation:

1. Is it possible to obtain polynomial<sup>3</sup> sample complexity for off-policy function estimation, using function classes

that only satisfy realizability-type assumptions?

Can we specify a distribution ν to the estimation algorithm, such that the learned function will be close to the groundtruth under ν?

In this work, we answer both open questions in the positive. By imposing proper regularization on the MIS objectives, we provide off-policy function estimation guarantees under only realizability assumptions on the function classes. Compared to commonly used regularization in MIS (Nachum et al., 2019a; Nachum & Dai, 2020; Yang et al., 2020), our regularizer is much more flexible and can account for an arbitrary user-specified distribution  $\nu$ , under which the learned function will be close to the groundtruth. We provide exact characterization of the optimal dual solution that needs to be realized by the discriminator, which determines the datacoverage assumption in value-function learning. As another surprising observation, the regularizer can be altered to relax the data-coverage requirement, and in the ideal case completely eliminate it when strong side information is available. Proof-of-concept experiments are also conducted to validate our theoretical predictions.

# 2. Related Works

**Regularization in MIS** The use of regularization is very common in the MIS literature, especially in DICE algorithms (Nachum et al., 2019a;b; Yang et al., 2020). However, most prior works that consider regularization use tabular derivations and seldom provide finite-sample functionapproximation guarantees on even return estimation, let alone function estimation. (An exception is the work of Uehara et al. (2021), who analyze related estimators under Bellman-completeness-type assumptions; see the next paragraph.) More importantly, while some existing DICE estimators are subsumed as our special cases when we choose very simple regularizers (see Remark 3 in Section 5), prior works provide very limited understanding in how the choice of regularizers affects learning guarantees, and hence have only considered these naïve forms of regularization (typically state-action-independent and under  $d^D$ )—as different forms of regularization are essentially treated equally under a coarse-grained theory (Yang et al., 2020). In contrast, we provide much more fine-grained characterization of the effects of regularization, which leads to novel insights about how to design better regularizers.

**Fitted-Q Evaluation (FQE)** Outside the MIS literature, one can obtain return *and* value-function estimation guarantees via FQE (Duan et al., 2020; Chen & Jiang, 2019; Le et al., 2019; Uehara et al., 2021). However, it is well understood that FQE and related approaches require Bellmancompleteness-type assumptions, such as the function class being *closed* under the Bellman operator. Even putting aside the difference between completeness vs. realizability, we al-

<sup>&</sup>lt;sup>2</sup> While the occupancies of these policies may not be known in general, they may be estimated from samples or approximated using domain knowledge. We give a more nuanced discussion in the conclusion.

<sup>&</sup>lt;sup>3</sup> By "polynomial", we mean polynomial in the horizon, the statistical capacities and the boundedness of the function classes, and the parameter that measures the degree of data coverage.

low for a user-specified error-measuring distribution, which is not available in FQE or any other existing method. The only distribution these methods are aware of is the data distribution  $d^D$ , and even so, FQE and variants rarely provide guarantees on  $\|\hat{q} - q^{\pi}\|_{2,d^D}$ , but often on the Bellman error (e.g.,  $\|\hat{q} - \mathcal{T}^{\pi}\hat{q}\|_{2,d^D}$ ) instead (Uehara et al., 2021), and obtaining guarantees on a distribution of interest often requires multiple indirect translations and loose relaxations.

**LSTDQ** Our analyses focus on general function approximation. When restricted to linear classes, function estimation guarantees for  $q^{\pi}$  under  $d^{D}$  can be obtained by LSTDQ methods (Lagoudakis & Parr, 2003; Bertsekas & Yu, 2009; Dann et al., 2014) when the function class only satisfies realizability of  $q^{\pi}$  (Perdomo et al., 2022). However, this requires an additional matrix invertibility condition (see Assumption 3 of Perdomo et al. (2022)), and it is still unclear what this condition corresponds to in general function approximation <sup>4</sup>. Moreover, many general methods—including MIS (Uehara et al., 2020) and other minimax methods (Antos et al., 2008; Xie et al., 2021)—coincide with LSTDQ in the linear case, so the aforementioned results can be viewed as a specialized analysis leveraging the properties of linear classes.

**PRO-RL (Zhan et al., 2022)** Our key proof techniques are adapted from Zhan et al. (2022), whose goal is offline policy learning. They learn the importance weight function  $w^{\pi}$  for a near-optimal  $\pi$ , and provide  $\|\widehat{w} - w^{\pi}\|_{2,d^{D}}$  guarantees as an intermediate result. Despite using similar technical tools, our most interesting and surprising results are in the value-function estimation setting, which is not considered by Zhan et al. (2022). Our novel algorithmic insights, such as incorporating error-measuring distributions and approximate models in the regularizers, are also potentially useful in Zhan et al. (2022)'s policy learning setting. Our analyses also reveal a number of important differences between OPE and offline policy learning, which will be discussed in Appendix A.

#### 3. Preliminaries

We consider off-policy evaluation (OPE) in Markov Decision Processes (MDPs). An MDP is specified by its state space S, action space A, transition dynamics  $P: S \times A \rightarrow \Delta(S)$  ( $\Delta(\cdot)$  is the probability simplex), reward function  $R: S \times A \rightarrow \Delta([0, 1])$ , discount factor  $\gamma \in [0, 1)$ , and an initial state distribution  $\mu_0 \in \Delta(S)$ . We assume S and A are finite and discrete, but their cardinalities can be arbitrarily large. Given a target policy  $\pi: S \rightarrow \Delta(A)$ , a random trajectory  $s_0, a_0, r_0, s_1, a_1, r_1, \ldots$  can be generated as  $s_0 \sim \mu_0, a_t \sim \pi(\cdot|s_t), r_t \sim R(\cdot|s_t, a_t), s_{t+1} \sim P(\cdot|s_t, a_t)$ ,

 $\forall t \geq 0$ ; we use  $\mathbb{E}_{\pi}$  and  $\mathbb{P}_{\pi}$  to refer to expectation and probability under such a distribution. The expected discounted return (or simply return) of  $\pi$  is  $J(\pi) := \mathbb{E}_{\pi}[\sum_{t} \gamma^{t} r_{t}]$ . The Q-value function of  $\pi$  is the unique solution of the Bellman equations  $q^{\pi} = \mathcal{T}^{\pi}q^{\pi}$ , with the Bellman operator  $\mathcal{T}^{\pi}$ :  $\mathbb{R}^{S \times \mathcal{A}} \to \mathbb{R}^{S \times \mathcal{A}}$  defined as  $\forall q \in \mathbb{R}^{S \times \mathcal{A}}, (\mathcal{T}^{\pi}q)(s, a) :=$  $\mathbb{E}_{r \sim R(\cdot|s,a)}[r] + \gamma(P^{\pi}q)(s, a)$ . Here  $P^{\pi} \in \mathbb{R}^{|S \times \mathcal{A}| \times |S \times \mathcal{A}|}$ is the state-action transition operator  $\pi$ , defined as  $(P^{\pi}q)(s, a) := \mathbb{E}_{s' \sim P(\cdot|s,a),a' \sim \pi(\cdot|s')}[q(s', a')]$ . Functions over  $S \times \mathcal{A}$  (such as q) are also treated as  $|S \times \mathcal{A}|$ dimensional vectors interchangeably.

In OPE, we want to estimate  $q^{\pi}$  and other functions of interest based on a historical dataset collected by a possibly different policy. As a standard simplification, we assume that the offline dataset consisting of n i.i.d. tuples  $\{(s_i, a_i, r_i, s_i')\}_{i=1}^n$  sampled as  $(s_i, a_i) \sim d^D$ , r = $R(s_i,a_i)$ , and  $s_i' \sim P(\cdot|s_i,a_i)$ . We call  $d^D$  the (offline) data distribution. As another function of interest, the (marginalized importance) weight function  $w^{\pi}$  is defined as  $w^{\pi}(s, a) := d^{\pi}(s, a)/d^{D}(s, a)$ , where  $d^{\pi}(s, a) =$  $(1-\gamma)\sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{\pi}[s_t = s, a_t = a]$  is the discounted state-action occupancy of  $\pi$ . For technical convenience we assume  $d^D(s,a) > 0 \ \forall s,a$ , so that quantities like  $w^{\pi}$  are always well defined and finite.<sup>5</sup> Similarly to  $q^{\pi}$ ,  $w^{\pi}$  also satisfies a recursive equation, inherited from the Bellman flow equation for  $d^{\pi}$ :  $d^{\pi} = (1 - \gamma)\mu_0^{\pi} + \gamma \tilde{P}^{\pi} d^{\pi}$ , where  $(s,a) \sim \mu_0^{\pi} \Leftrightarrow s \sim \mu_0, a \sim \pi(\cdot|s)$  is the initial state-action distribution, and  $\widetilde{P}^{\pi} = (P^{\pi})^{\top}$  is the transpose of the transition matrix.

**Function Approximation** We will use function classes Q and W to approximate  $q^{\pi}$  and  $w^{\pi}$ , respectively. We assume finite Q and W, and extension to infinite classes under appropriate complexity measures (e.g., covering number) is routine and orthogonal to the main insights of the paper.

Additional Notation  $\|\cdot\|_{2,\nu} := \sqrt{\mathbb{E}_{\nu}[(\cdot)^2]}$  is the weighted 2-norm of a function under distribution  $\nu$ . We also use a standard shorthand  $f(s,\pi) := \mathbb{E}_{a \sim \pi(\cdot|s)}[f(s,a)]$ .  $u \circ v$  between two vectors u and v of the same dimension is elementwise multiplication, and u/v is elementwise division.

# 4. Value-function Estimation

In this section we show how to estimate  $\hat{q} \approx q^{\pi}$  with guarantees on  $\|\hat{q} - q^{\pi}\|_{2,\nu}$  for a user-specified  $\nu$ , and identify the assumptions under which provable sample-complexity guarantees can be obtained. We begin with the familiar

<sup>&</sup>lt;sup>4</sup> It is hinted by Uehara et al. (2020) that the invertibility is related to a loss minimization condition in MIS, but the connection only holds for return estimation.

<sup>&</sup>lt;sup>5</sup> It will be trivial to remove this assumption at the cost of cumbersome derivations. Also, these density ratios can still take prohibitively large values even if they are finite, and we will need to make additional boundedness assumptions to enable finite-sample guarantees anyway, so their finiteness does not trivialize the analyses.

Bellman equations, that  $q^{\pi}$  is the unique solution to:

$$\mathbb{E}_{r \sim R(\cdot|s,a)}[r] + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}[q(s',\pi)] - q(s,a) = 0,$$
  
$$\forall s, a \in \mathcal{S} \times \mathcal{A}.$$
  
(1)

While the above set of equations uniquely determines  $q = q^{\pi}$ , this is only true if we can enforce *all* the  $|S \times A|$  constraints, which is intractable in large state-space problems. In fact, even estimating (a candidate *q*'s violation of) a single constraint is infeasible as that requires sampling from the same state multiple times, which is related to the infamous double-sampling problem (Baird, 1995).

To overcome this challenge, prior MIS works often relax (1) by taking a *weighted combination* of these equations, e.g.

$$\mathbb{E}_{d^{D}}\left[w(s,a)\left(r(s,a) + \gamma q(s',\pi) - q(s,a)\right)\right] = 0,$$
  
$$\forall w \in \mathcal{W}. \quad (2)$$

In words, instead of enforcing  $|S \times A|$  equations, we only enforce their linear combinations; the linear coefficients are  $d^D(s, a) \cdot w(s, a)$ , and w belongs to a class W with limited statistical capacity to enable sample-efficient estimation. While each constraint in (2) now be efficiently checked on data, this comes with a big cost that a solution to (2) is not necessarily  $q^{\pi}$ . Prior works handle this dilemma by aiming lower: instead of learning  $\hat{q} \approx q^{\pi}$ , the problem becomes tractable if we only aim to learn  $\hat{q}$  that can approximate the policy's return, i.e.  $\mathbb{E}_{s \sim \mu_0}[\hat{q}(s,\pi)] \approx$  $J(\pi) = \mathbb{E}_{s \sim \mu_0}[q^{\pi}(s,\pi)]$ . While Uehara et al. (2020) show that this is possible under  $w^{\pi} \in \mathcal{W}$ , they also show explicit counterexamples where  $\hat{q} \neq q^{\pi}$  even with infinite data. As a result, how to estimate  $\hat{q} \approx q^{\pi}$  under comparable assumptions (instead of the prohibitive  $\mathcal{W} = \mathbb{R}^{S \times A}$  as in Proposition 1) is still an open problem.

#### 4.1. Estimator

We now describe our approach to this problem. Recall that the goal is to obtain error bounds for  $\|\widehat{q} - q^{\pi}\|_{2,\nu}$  for some distribution  $\nu \in \Delta(S \times A)$  provided by the user. Note that we do *not* require information about r and s' that are generated after  $(s, a) \sim \nu$  and only care about the (s, a)marginal itself, so the user can pick  $\nu$  without knowing the transition and the reward functions of the MDP. We assume that  $\nu$  is given in a way that we can take its expectation  $\mathbb{E}_{(s,a)\sim\nu}[(\cdot)]$ ; the extension to the case where  $\nu$  is given via samples is straightforward.

To achieve this goal, we first turn (1) into an equivalent *constrained convex program*: given a collection of strongly convex and differentiable functions  $f = \{f_{s,a} : \mathbb{R} \to \mathbb{R}\}_{s,a}$ 

we will later discuss its choice-consider

$$\min_{q} \mathbb{E}_{(s,a)\sim\nu}[f_{s,a}(q(s,a))] \tag{3}$$

s.t. 
$$\mathbb{E}_{r \sim R(\cdot|s,a)}[r] + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}[q(s',\pi)] - q(s,a) = 0,$$
  
 $\forall s, a.$ 

The constraints here are the same as (1). Since (1) uniquely determines  $q = q^{\pi}$ , the feasible space of (3) is a *singleton*, so we can impose any objective function on top of these constraints, and it will not change the optimal solution (which is always  $q^{\pi}$ , the only feasible point). Here, we use an objective function of the form  $\mathbb{E}_{(s,a)\sim\nu}[f_{s,a}(q(s,a))]$ , where  $\nu$  is the user-provided distribution of interest. Choosing  $f_{s,a}(q(s,a)) = 0 \forall s, a$ , for example, recovers MQL in Uehara et al. (2020). With this choice, however, as shown in Proposition 1,  $q^{\pi}$  cannot be identified without prohibitive assumptions. As we will see, other choices of f will serve as important regularizers in the function-approximation setting, and will be crucial for our function estimation guarantees.

**Remark 1** ((s, a)-dependence of f). Regularizers in prior works are (s, a)-*independent* (Nachum et al., 2019a; Yang et al., 2020; Zhan et al., 2022). As we will see in Section 4.3, allowing for (s, a)-dependence is very important for designing regularizers with improved guarantees and performances.

We now rewrite (3) in its Lagrangian form, with  $d^D \circ w$  serving the role of dual variables:

$$\min_{q} \max_{w} L_{f}^{q}(q, w) := \mathbb{E}_{\nu}[f_{s,a}(q(s, a))] + \mathbb{E}_{d^{D}} [w(s, a) (r(s, a) + \gamma q(s', \pi) - q(s, a))].$$
(4)

Finally, our actual estimator approximates (4) via finitesample approximation of the population loss  $L_f^q$ , and searches over restricted function classes Q and W for qand w, respectively:

$$(\widehat{q}, \widehat{w}) = \operatorname*{arg\,min}_{q \in \mathcal{Q}} \operatorname*{arg\,max}_{w \in \mathcal{W}} \widehat{L}_{f}^{q}(q, w), \tag{5}$$

where 
$$L_f^q(q, w) := \mathbb{E}_{\nu}[f_{s,a}(q(s, a))]$$
  
+  $\frac{1}{n} \sum_{i=1}^n w(s_i, a_i) (r_i + \gamma q(s'_i, \pi) - q(s_i, a_i))$ 

**Intuition for identification** Before giving the detailed finite-sample analysis, we provide some high-level intuitions for why we can obtain the desired guarantee on  $\|\hat{q} - q^{\pi}\|_{2,\nu}$ . Note that (4) is structurally similar to (2), and we still cannot verify the Bellman equation for  $q^{\pi}$  in a per-state-action manner, so the caveat of (2) seems to remain; why can we identify  $q^{\pi}$  under  $\nu$ ?

The key here is to show that it suffices to check the loss function  $L_f^q$  only under a special choice of w (as opposed to all of  $\mathbb{R}^{S \times \mathcal{A}}$ ). Importantly, this special w is *not*  $w = w^{\pi}$ ; rather, it is the function  $w_f^*$  in the saddle-point solution of our regularized objective,  $(q^{\pi}, w_f^*) = \arg \min_q \arg \max_w L_f^q(q, w)$ , for which we later provide the closed form. As long as  $w_f^* \in \mathcal{W}$ —even if  $\mathcal{W}$  is extremely "simple" and contains nothing but  $w_f^*$ —we can identify  $q^{\pi}$ . Intuitively,  $w^{\pi}$  should not appear in our analysis at all:  $w^{\pi}$  is defined with respect to the initial distribution of the MDP  $\mu_0$ , which has nothing to do with our goal of bounding  $\|\hat{q} - q^{\pi}\|_{2,\nu}$ .

To see that, it is instructive to consider the special case of  $\mathcal{W} = \{w_f^*\}$  and the limit of infinite data. In this case, our estimator becomes  $\arg \min_{q \in \mathcal{Q}} L_f^q(q, w_f^*)$ . By the definition of saddle point:

$$L_f^q(q^{\pi}, w_f^*) \leq L_f^q(q, w_f^*), \ \forall q$$

While this shows that  $q^{\pi}$  is a minimizer of the loss, it does not imply that it is a unique minimizer. However, identification immediately follows from the convexity brought by regularization: since  $f : \mathbb{R} \to \mathbb{R}$  is strongly convex,  $q \to \mathbb{E}_{\nu}[f_{s,a}(q(s, a))]$  as a mapping from  $\mathbb{R}^{S \times A}$  to  $\mathbb{R}$  is strongly convex under  $\|\cdot\|_{2,\nu}$  (see Lemma 7 in Appendix B for a formal statement and proof), and  $L_{f}^{q}(q, w_{f}^{*})$  inherits such convexity since the other terms are affine in q. It is then obvious that  $q^{\pi}$  is the unique minimizer of  $L_{f}^{q}(q^{\pi}, w_{f}^{*})$  up to  $\|\cdot\|_{2,\nu}$ , that is, any minimizer of  $L_{f}^{q}$  must agree with  $q^{\pi}$  on (s, a) pairs supported on  $\nu$ . Our finite-sample analysis below shows that the above reasoning is robust to finite-sample errors and having functions other than  $w_{f}^{*}$  in  $\mathcal{W}$ .

#### 4.2. Finite-sample Guarantees

In this subsection we state the formal guarantee of our estimator for  $q^{\pi}$ , and the assumptions under which it holds. We start with the condition on the regularization function f:

Assumption 1 (Strong convexity of f). Assume  $f_{s,a} : \mathbb{R} \to \mathbb{R}$  is nonnegative, differentiable, and  $M^q$ -strongly convex for each  $s \in S, a \in A$ . In addition, assume both  $f_{s,a}$  and  $f'_{s,a}$  take finite values for any finite input.

This assumption can be satisfied by a simple choice of  $f_{s,a}(x) = \frac{1}{2}x^2$ , which is independent of (s, a) and yields  $M^q = 1$ . Alternative choices of f are discussed in Section 4.3. Next are the realizability and boundedness of W and Q:

Assumption 2 (Realizability). Suppose  $w_f^* \in \mathcal{W}, q^{\pi} \in \mathcal{Q}$ .

**Assumption 3** (Boundedness of W and Q). Suppose W and Q are bounded, that is,

$$C_{\mathcal{Q}}^{q} := \max_{q \in \mathcal{Q}} \|q\|_{\infty} < \infty,$$
  
$$C_{\mathcal{W}}^{q} := \max_{w \in \mathcal{W}} \|w\|_{\infty} < \infty.$$

As a remark, Assumption 2 implicitly assumes the existence of  $w_t^*$ . As we will see in Section 4.3, the existence

and finiteness of  $w_f^*$  is automatically guaranteed given the finiteness of  $f'_{s,a}$  (Assumption 1) and  $d^D(s,a) > 0 \ \forall s,a$ . More importantly, Assumptions 2 and 3 together imply that  $\|q^{\pi}\|_{\infty} \leq C_Q^q$  and  $\|w_f^*\|_{\infty} \leq C_W^q$ , which puts constraints on how small  $C_Q^q$  and  $C_W^q$  can be. For example, it is common to assume that  $C_Q^q = \frac{1}{1-\gamma}$ , i.e., the maximum possible return when rewards are bounded in [0, 1], and this way  $\|q^{\pi}\|_{\infty} \leq C_Q^q$  will hold automatically. The magnitude of  $\|w_f^*\|_{\infty}$  and  $C_W^q$ , however, is more nuanced and interesting, and we defer the discussion to Section 4.3.

Now we are ready to state the main guarantee for identifying  $q^{\pi}$ . All proofs of this section can be found in Appendix B.

**Theorem 2.** Suppose Assumptions 1, 2, 3 hold. Then, with probability at least  $1 - \delta$ ,

$$||\widehat{q} - q^{\pi}||_{2,\nu} \le 2\sqrt{\frac{\epsilon_{stat}^q}{M^q}},$$

where 
$$\epsilon_{stat}^q = \left(C_{\mathcal{W}}^q + (1+\gamma)C_{\mathcal{W}}^q C_{\mathcal{Q}}^q\right) \sqrt{2\log \frac{2|\mathcal{W}||\mathcal{Q}|}{\delta}/n}.$$

Theorem 2 shows the desired bound on  $\|\hat{q} - q^{\pi}\|_{2,\nu}$ , which depends on the magnitude of functions in  $\mathcal{W}$  and  $\mathcal{Q}$  as well as their logarithmic cardinalities, which are standard measures of statistical complexity for finite classes. One notable weakness is the  $O(n^{-1/4})$  slow rate; this is due to translating the  $\epsilon_{stat}^q = O(n^{-1/2})$  deviation between L and  $\hat{L}$  into  $\|\hat{q} - q^{\pi}\|_{2,\nu}$  via a convexity argument, which takes square root of the error. The possibility of and obstacles to obtaining an  $O(n^{-1/2})$  rate will be discussed in Section 7.

# **4.3.** On the Closed Form of $w_f^*$ and the Data Coverage Assumptions

One unusual aspect of our guarantees in Section 4.2 is that we do not make any explicit data coverage assumptions, yet such assumptions are known to be necessary even for return estimation (typically the boundedness of  $w^{\pi} = d^{\pi}/d^D$ ). Indeed, our data-coverage assumption is implicit in Assumptions 2 and 3, which require  $||w_f^*||_{\infty} \leq C_W^q < \infty$ . If data fails to provide sufficient coverage,  $||w_f^*||_{\infty}$  will be large and our bound in Theorem 2 will suffer due to a large value of  $C_W^q$ .

To make the data coverage assumption explicit, we provide the closed-form expression of  $w_f^*$ :

**Lemma 3.** The saddle point of (4) is  $(q^{\pi}, w_f^*) = \arg \min_q \arg \max_w L_f^q(q, w)$ , where

$$w_f^* = (I - \gamma \widetilde{P}^{\pi})^{-1} \left( \nu \circ f'(q^{\pi}) \right) / d^D.$$
 (6)

Here  $f'(q^{\pi})$  is the shorthand for  $[f'_{s,a}(q^{\pi}(s,a))]_{s,a} \in \mathbb{R}^{S \times A}$ .

The closed-form expression in (6) looks very much like a density ratio: if we replace  $\nu \circ f'(q^{\pi})$  with  $\mu_0^{\pi}$ , we have

 $(I - \gamma \widetilde{P}^{\pi})^{-1} \mu_0^{\pi} = d^{\pi}/(1 - \gamma)$ , and the expression would be the ratio between  $d^{\pi}$  and  $d^D$  (up to a horizon factor). Therefore,  $w_f^*$  can be viewed as the density ratio of  $\pi$  against  $d^D$ when  $\pi$  starts from the "fake" initial distribution  $\nu \circ f'(q^{\pi})$ . However,  $\nu \circ f'(q^{\pi})$  is in general not a valid distribution, as it is not necessarily normalized or even non-negative, making  $||w_f^*||_{\infty}$ , which are more interpretable and give novel insights into how to relax the data-coverage assumption via tweaking f.

**Proposition 4.**  $||w_f^*||_{\infty} \leq \frac{1}{1-\gamma} \cdot ||d_{\nu}^{\pi}/d^D||_{\infty} \cdot ||f'(q^{\pi})||_{\infty}$ , where  $d_{\nu}^{\pi}$  is the discounted state-action occupancy of  $\pi$  under  $\nu$  as the initial state-action distribution.

The proposition states that  $||w_f^*||_{\infty}$  can be bounded if data provides sufficient coverage over  $d_{\nu}^{\pi}$ , and if  $f'(q^{\pi})$  is bounded. The former shows that  $d^D$  needs to cover not only  $\nu$ , but also state-action pairs reachable by  $\pi$  starting from  $\nu$ . The latter is easily satisfied, and can be bounded again for concrete choices of f, e.g.  $||f'(q^{\pi})||_{\infty} \leq ||q^{\pi}||_{\infty} \leq \frac{1}{1-\gamma}$ for  $f_{s,a}(x) = \frac{1}{2}x^2$ .

**Designing** f to relax the coverage assumption Lemma 3 shows that the coverage assumption  $(||w_f^*||_{\infty})$  depends on f (or rather its derivative f'), which opens up the possibility of properly designing f to relax it. In fact, we could completely eliminate the coverage assumption if we could set  $f'(q^{\pi}) = \mathbf{0}$ , but that would require unrealistically strong side information.

As a concrete example, consider  $f_{s,a}(x) = \frac{1}{2}(x - q^{\pi}(s,a))^2$ , and it is easy to verify that  $f'_{s,a}(q^{\pi}(s,a)) = x - q^{\pi}(s,a)|_{x=q^{\pi}(s,a)} = 0$ . Compared to  $f_{s,a}(x) = \frac{1}{2}x^2$ , the new f essentially adds a 1st-order term  $q^{\pi}(s,a) \cdot x$  to change  $w_f^*$ , while leaving the convexity required by Assumption 1 intact, which only depends on the 2nd-order term  $\frac{1}{2}x^2$ . Of course, this is not a viable choice of f in practice as it requires knowledge of  $q^{\pi}$ , which is precisely our learning target.

While the reason  $f_{s,a}(x) = \frac{1}{2}(x - q^{\pi}(s, a))^2$  can eliminate the coverage requirement is obvious retrospectively  $(q^{\pi} \text{ already minimizes } \mathbb{E}_{\nu}[f(q)]$  even without any data), our analyses apply much more generally and characterize the effects of arbitrary f on the coverage assumption. Inspired by this example, we can consider practically feasible choices such as  $f_{s,a}(x) = \frac{1}{2}(x - \tilde{q}(s, a))^2$ , where  $\tilde{q}$  is an approximation of  $q^{\pi}$  obtained by other means, e.g. a guess based on domain knowledge. If  $\tilde{q} \approx q^{\pi}$ , our estimator enjoys significantly relaxed coverage requirements. But even if  $\tilde{q}$  is a poor approximation of  $q^{\pi}$ , it does not affect our estimation guarantees as long as the condition implied by Proposition 4 is satisfied. (In fact,  $f_{s,a}(x) = \frac{1}{2}x^2$  is a special case of  $\tilde{q} \equiv 0$ .) Such a use of approximate models is similar to how doubly robust estimators (Dudík et al., 2011; Jiang & Li,

2016; Thomas & Brunskill, 2016) enjoy reduced variance given an accurate model, and remain unbiased even if the approximate model is arbitrarily poor. We will also empirically evaluate the effectiveness of this idea in Section 6.

#### 5. Weight-function Estimation

Similar to value-function estimation, our methodology can also be applied to estimate the weight function  $w^{\pi}$ . Due to the similarity with Section 4 in the high-level spirit, we will be concise in this section and only explain in detail when there is a conceptual difference from Section 4. Some notations (such as the function classes W and Q) will be abused, but we emphasize that this section considers a different learning task than Section 4, so they should be viewed as different objects (e.g., the realizability assumptions for W and Q below will be different from those in Section 4).

As before, we assume that the user provides a distribution<sup>6</sup>  $\eta \in \Delta(S \times A)$  and our goal is to develop an estimator with guarantees on  $\|\widehat{w} - w^{\pi}\|_{2,\eta}$ . Analogous to Section 4, consider

$$\min_{w} \mathbb{E}_{(s,a)\sim\eta}[f_{s,a}(w(s,a))]$$
s.t.  $d^{D}(s,a)w(s,a) = (1-\gamma)\mu_{0}^{\pi}(s,a)$   
 $+\gamma \sum_{s',a'} P^{\pi}(s,a|s',a')d^{D}(s',a')w(s,a),$   
 $\forall s, a.$ 
(7)

Here  $f = \{f_{s,a}\}_{s,a}$  will need to satisfy similar assumptions as in Section 4. The constraints are the Bellman flow equations with a change of variable  $d(s,a) = d^D(s,a) \cdot w(s,a)$ . Their unique solution is  $d(s,a) = d^{\pi}(s,a)$  (and hence  $w(s,a) = d^{\pi}(s,a)/d^D(s,a)$ ), thus the feasible space is again a singleton, and the objective does not alter the optimal solution. We then use dual variables q to rewrite (7) in its Lagrangian form:

$$\min_{w} \max_{q} L_{f}^{w}(q, w)$$

$$:= \mathbb{E}_{\eta}[f_{s,a}(w(s, a))] + (1 - \gamma)\mathbb{E}_{\mu_{0}}[q(s, \pi)]$$

$$+ \mathbb{E}_{d^{D}}[w(s, a)(\gamma q(s', \pi) - q(s, a))]$$
(8)

We approximate the saddle-point solutions by optimizing the empirical loss  $\hat{L}_{f}^{w}$  over restricted function classes  $\mathcal{W}, \mathcal{Q}$ :

$$(\widehat{w}, \widehat{q}) = \operatorname*{arg\,max}_{w \in \mathcal{W}} \operatorname*{arg\,min}_{q \in \mathcal{Q}} \widehat{L}_{f}^{w}(q, w),$$

where  $\widehat{L}_{f}^{w}(q, w) := \mathbb{E}_{\eta}[f_{s,a}(w(s, a))] + \frac{1-\gamma}{n_{0}} \sum_{j=1}^{n_{0}} q(s_{j}, \pi) + \frac{1}{n} \sum_{i=1}^{n} w(s_{i}, a_{i}) (\gamma q(s'_{i}, \pi) - q(s_{i}, a_{i})),$ 

<sup>&</sup>lt;sup>6</sup> Recall we assume  $d^{D}(s, a) > 0 \ \forall s, a$  for technical convenience. When this is not the case,  $\eta$  should be supported on  $d^{D}$ , as the target function  $w^{\pi}$  is only defined on the support of  $d^{D}$ .

v

and  $\{s_j\}_{j=1}^{n_0}$  is a separate dataset sampled i.i.d. from  $\mu_0$  to provide information about initial distribution.

We provide the closed-form expression for the saddle point of  $L_f^w$  below, which resembles the Q-function for a proxy reward function  $f'(w^{\pi}) \circ \eta/d^D$ .

**Lemma 5.** The closed form solutions of (8) are  $(w^{\pi}, q_{f}^{*}) =$  $\arg\min_{w} \arg\max_{q} L_{f}^{w}(q, w)$ , where

$$q_f^* = (I - \gamma P^{\pi})^{-1} (f'(w^{\pi}) \circ \eta / d^D).$$
(9)

Remark 2 (Data Coverage Assumption). As we will see, the only data coverage assumption we need is the boundedness of  $w^{\pi} = d^{\pi}/d^{\bar{D}}$ . Since  $w^{\pi}$  is the function of interest and practical algorithms can only output functions of wellbounded ranges, such an assumption is an essential part of the learning task itself and hardly an additional requirement. Moreover, unlike Section 4, changing f here will not affect the data-coverage assumption, though it still alters  $q_f^*$ , and a properly chosen f (e.g., with  $f'(w^{\pi}) \approx 0$ ) can still result in a  $q_f^*$  with small magnitude and thus make learning easier.

Remark 3 (Connection to DualDICE). We can recover DualDICE (Nachum et al., 2019a) by choosing  $f_{s,a}(x) =$  $\frac{1}{2}x^2$  and  $\nu = d^D$ . Despite producing the same estimator, the derivations and assumptions under which the two works analyze the estimator are different. Their Theorem 2 only provides return estimation guarantees, and depends on an implicit assumption of highly expressive function classes<sup>7</sup> similar to Proposition 1. Moreover, they do not characterize how the choice of f can affect the learning guarantees (their f is (s, a)-independent). This is one of the main insights of our paper and leads to the discovery of more practical regularizers, e.g.  $f_{s,a}(x) = \frac{1}{2}(x - \tilde{w}(s, a))^2$  with model  $\tilde{w}$ .

Below we present the assumptions, then learning guarantee for  $\widehat{w}$ .

Assumption 4 (Strongly Convex Objective). Suppose for all s, a,  $f_{s,a}$  is differentiable, non-negative, and  $M^w$ strongly convex. Further, suppose  $f_{s,a}$  and its derivative take finite values on any finite inputs, and let  $C_f^w :=$  $\max_{w \in \mathcal{W}} ||f(w)||_{\infty}.$ 

Assumption 5 (Realizability). Suppose  $w^{\pi} \in \mathcal{W}, q_f^* \in \mathcal{Q}$ .

Assumption 6 (Bounded W and Q). Let  $C_W^w$  :=  $\max_{w \in \mathcal{W}} ||w||_{\infty}$  and  $C_{\mathcal{Q}}^w := \max_{q \in \mathcal{Q}} ||q||_{\infty}$ . Suppose  $\mathcal{W}$ and Q are bounded function classes, that is,  $C_{\mathcal{W}}^w < \infty$  and  $C_{\mathcal{Q}}^w < \infty.$ 

**Theorem 6.** Suppose Assumptions 4, 5, 6, hold. Then w.p.  $> 1 - \delta$ .

$$\|\widehat{w} - w^{\pi}\|_{2,\eta} \le 2\sqrt{\frac{\epsilon_{stat}^w}{M^w}},$$

where 
$$\epsilon_{stat}^w = (1+\gamma)C_{\mathcal{W}}^w C_{\mathcal{Q}}^w \sqrt{2\log\frac{4|\mathcal{Q}||\mathcal{W}|}{\delta}}/n + (1-\gamma)C_{\mathcal{Q}}^w \sqrt{2\log\frac{4|\mathcal{Q}|}{\delta}}/n_0.$$

# 6. Experiments

We now provide experimental results to verify our theoretical predictions and insights. As Yang et al. (2020) have performed extensive experiments on return estimation with simple regularization  $(f_{s,a}(x) = \frac{1}{2}x^2)$ , we focus on the task of  $q^{\pi}$  estimation, and the following two questions unique to our work:

**Q1.** When the goal is to minimize  $\|\hat{q} - q^{\pi}\|_{2,\nu}$ , how much benefit does regularizing with  $\nu$  bring in practice, compared to regularizing with other distributions (or no regularization at all)?

Q2. Can incorporating (even relatively poor) models in regularization (e.g.,  $f_{s,a}(x) = \frac{1}{2}(x - \tilde{q}(s,a))^2$  from Section 4.3) improve estimation?

Setup We study these questions in a large tabular Gridwalk environment (Nachum et al., 2019a; Yang et al., 2020), with a deterministic target policy  $\pi$  that is optimal, and a behavior policy that provides limited coverage over the target policy; see Appendix D for details. To mimic the identification challenges associated with restricted function classes, we use a linear function class  $Q = \{ \Phi^\top \alpha : \alpha \in \mathbb{R}^d \}$  and discriminator class  $\mathcal{W} = \{ \widetilde{\Phi}^\top \beta : \beta \in \mathbb{R}^k \}$ , where  $k < d \ll |\mathcal{S} \times \mathcal{A}|$ . The features  $\Phi \in \mathbb{R}^{|S \times A| \times d}$ ,  $\widetilde{\Phi} \in \mathbb{R}^{|S \times A| \times k}$  are chosen to satisfy the realizability assumptions of all estimators. Under linear classes, our estimator ((5)) becomes a convex optimization problem with d variables and k linear constraints, and can be solved by standard packages. This allows us to avoid difficult minimax optimization-which is still an open problem in the MIS literature-and focus on the statistical behaviors of our estimators, which is what our theoretical predictions are about.

Remark 4. When no regularization is used, our estimator coincides with MOL (Uehara et al., 2020). If we further had  $\Phi = \Phi$ , the estimator would coincide with LSTDQ. While Section 2 mentioned that LSTDQ enjoys functionestimation guarantees (Perdomo et al., 2022) (and folklore suggests they extend to  $\Phi \neq \Phi$ ), the guarantee only holds in the regime of  $k \ge d$ , i.e., the k linear constraints are over-determined. In our case, however, we have underdetermined constraints (k < d), creating a more challenging learning task (which our theory can handle) where LSTDQ's guarantees do not apply.

Choice of Distributions We consider a set of diverse distri-

<sup>&</sup>lt;sup>7</sup> In our notation, they measure the approximation error of  $\mathcal{W}$  as  $\max_{w' \in \mathbb{R}^{S \times A}} \min_{w \in \mathcal{W}} \|w - w'\|$ , essentially requiring  $\mathcal{W}$  (and similarly Q) to closely approximate every function over  $S \times A$ . However, we suspect that they could have measured realizability errors instead without changing much of their proofs.

**Beyond the Return: Off-policy Function Estimation** 



Figure 1. Error of off-policy return and function estimation as a function of sample size. Legend shows regularizing distribution  $\nu$  and header shows error-measuring distribution  $\nu'$  (see text). Error bars show 95% confidence intervals calculated from 1000 runs.



Figure 2. Estimation error when the regularizer incorporates a model  $\tilde{q}$ , where x-axes represent the parameter m that controls the quality of  $\tilde{q}$ . Sample size is 500 and the results are from 500 runs.

butions  $\mathcal{V} = \{d^D, \mu_0^{\pi}, d^{\pi}, U, p\}$ , where U is uniform over  $\mathcal{S} \times \mathcal{A}$  and  $p \propto (d^{\pi} \circ \mathbb{I}[w^{\pi} > 50])$ . The distribution p isolates the least-covered states reached by  $\pi$ , which makes learning an accurate Q-function on  $\nu$  a harder task.

**Results for Q1** We use a default regularizer  $f = \frac{1}{2}x^2$ with different regularizing distributions  $\nu \in \mathcal{V}$ , and measure  $\|\widehat{q} - q^{\pi}\|_{2,\nu'}$  for different  $\nu'$ . The results are shown in Figure 1, and exhibit the expected trend. For example, regularizing with  $\nu = p$  performs poorly when the error is measured under  $\nu' = d^{\overline{D}}$  and U due to the large mismatch between  $\nu$  and  $\nu'$ . However, when  $\nu' = p$  (rightmost panel), regularizing with  $\nu = p$  significantly outperforms others. Similar behaviors can also be observed on U, though they are certainly not absolute (e.g.,  $\nu = d^D$  does not do very well on  $\nu' = d^D$ ), which suggests potential directions for more refined theory. Moreover, using no regularization ("none") generally does not perform well for any  $\nu'$ , but still manages to achieve a high accuracy for return estimation  $J(\pi)$ , which is consistent with prior theory (Uehara et al., 2020) that return estimation does not require regularization.

**Results for Q2** We now use  $f_{s,a}(x) = \frac{1}{2}(x - \tilde{q}(s, a))^2$  with different  $\tilde{q}$  to verify how the quality of  $\tilde{q}$  affects estimation accuracy. We first consider a "uniform model"  $\tilde{q} = mq^{\pi} + (1 - m)\bar{q}$ , where  $\bar{q}$  is a constant and  $m \in [0, 1]$  controls the quality  $\tilde{q}$ . As shown in Panels 1 & 3 of Figure 2, our

estimator's accuracy generally improves with a better  $\tilde{q}$  (i.e., as *m* increases). Moreover, equipping  $\tilde{q}$  with an appropriate regularizing distribution  $\nu$  (e.g.,  $\nu = U$  for both panels) can significantly outperform no regularization, even with a very poor  $\tilde{q}$  (e.g., m = 0.1). It also outperforms the model prediction itself (i.e.,  $\hat{q} = \tilde{q}$ ), showing that the improvement is not from simply taking predictions from  $\hat{q}$ , but instead from using the regularization to better identify  $q^{\pi}$  from data.

The previous model's quality is uniform across  $S \times A$ . We then consider a scenario where  $\tilde{q}$  is zeroed out outside p's support, making it only a good approximation of  $q^{\pi}$  on p. In this case, we see that regularization cannot benefit much from the model when the error is measured on  $\nu' = U$  (Panel 2), but when  $\nu' = \nu = p$  (Panel 4), regularization can still bring benefits, as expected from our theory.

# 7. Discussion and Conclusion

In this paper, we showed that proper regularization can yield function-estimation guarantees for MIS methods under only realizable function approximation. Our results identify when and why accurate off-policy estimates of weight and value functions, which play important roles in larger reinforcement learning algorithms, can be obtained, and provide insight into how different regularizers lead to better estimates. Compared to prior works, our regularizer is more flexible and can accommodate a user-specified errormeasuring distribution. Further theoretical investigation provides fine-grained characterization of how the choice of regularization affects learning guarantees, which leads to the discovery of regularizers that incorporate approximate models (such as  $\tilde{q}$ ). While the superiority of such regularizers is perhaps obvious retrospectively, it is not allowed in the prior works' derivation that assumes (s, a)-independent regularization, and our theoretical results provide a deep understanding for even more general regularization schemes.

In Appendix A, we provide further discussions on two topics: (1) the barriers to obtaining a faster  $O(n^{-1/2})$  rate, and (2) comparison to Zhan et al. (2022) reveals interesting differences between off-policy function estimation and policy learning, and insights in this paper may also be useful for the policy learning task. In Appendix E we discuss robustness to approximation and optimization errors. Finally, in Appendix F, we discuss how estimated functions can be used in downstream off-policy evaluation tasks, and provide corresponding estimation guarantees.

For the aforementioned off-policy evaluation and off-policy learning (Liu et al., 2019) tasks, function estimation is (naturally) generally required to be accurate on  $\nu = d^D$  or  $\mu^{\pi}$ and  $\eta = d^D$  (Liu et al., 2018; Perdomo et al., 2022; Liu et al., 2019), and samples or exact distributions are accessible to the user. To this end, our method provides a concise answer to a missing piece of off-policy algorithms.

However, the interaction of our method with online training methods is less clear. As mentioned in the introduction, online algorithms using off-policy function estimation as a subroutine, such as (Kakade & Langford, 2002; Abbasi-Yadkori et al., 2019), may require the estimates to be accurate on unknown distributions such as  $d^{\pi}$  or  $d^{\pi^*}$  (where  $\pi^*$  is the optimal policy), which may not be immediately accessible to the user. If samples from a replay buffer are used to estimate functions of  $\pi$ , our method would provide estimation guarantees over  $\nu$  corresponding to the distribution of the buffer. Under appropriate coverage assumptions, a change of distribution could then be used to convert these guarantees over  $\nu$  to a guarantee over  $d^{\pi}$ . In other cases, the user may be able to use domain knowledge to "guess" a distribution  $\nu$  close to or covering the unknown distribution of interest. Then our guarantees for function estimation over  $\nu$  could similarly, with a change of distribution, be converted to guarantees on the true distribution of interest. To this end, an important avenue of future work involves a thorough investigation of how our off-policy function estimation method interacts with such online learning methods, their assumptions, and their guarantees.

#### Acknowledgment

The authors thank Wenhao Zhan and Jason Lee for valuable discussions during the early phase of the project. NJ acknowledges funding support from ARL Cooperative Agreement W911NF-17-2-0196, NSF IIS-2112471, NSF CA-REER IIS-2141781, and Adobe Data Science Research Award.

# References

- Abbasi-Yadkori, Y., Bartlett, P., Bhatia, K., Lazic, N., Szepesvari, C., and Weisz, G. Politex: Regret bounds for policy iteration using expert prediction. In *International Conference on Machine Learning*, pp. 3692–3702. PMLR, 2019.
- Agrawal, A., Verschueren, R., Diamond, S., and Boyd, S. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.
- Antos, A., Szepesvári, C., and Munos, R. Learning nearoptimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 2008.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceed*ings 1995, pp. 30–37. Elsevier, 1995.
- Bertsekas, D. P. and Yu, H. Projected equation methods for approximate solution of large linear systems. *Journal of Computational and Applied Mathematics*, 227(1):27–50, 2009.
- Chen, J. and Jiang, N. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, pp. 1042–1051. PMLR, 2019.
- Dann, C., Neumann, G., and Peters, J. Policy evaluation with temporal differences: A survey and comparison. *The Journal of Machine Learning Research*, 15:809–883, 2014.
- Diamond, S. and Boyd, S. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- Duan, Y., Jia, Z., and Wang, M. Minimax-optimal offpolicy evaluation with linear function approximation. In *International Conference on Machine Learning*, pp. 2701– 2709. PMLR, 2020.
- Dudík, M., Langford, J., and Li, L. Doubly Robust Policy Evaluation and Learning. In *Proceedings of the 28th International Conference on Machine Learning*, pp. 1097– 1104, 2011.

- Ernst, D., Geurts, P., and Wehenkel, L. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6:503–556, 2005.
- Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 652–661. PMLR, 2016.
- Kakade, S. and Langford, J. Approximately Optimal Approximate Reinforcement Learning. In *Proceedings of the 19th International Conference on Machine Learning*, volume 2, pp. 267–274, 2002.
- Kallus, N. and Uehara, M. Double reinforcement learning for efficient off-policy evaluation in markov decision processes. J. Mach. Learn. Res., 21(167):1–63, 2020.
- Lagoudakis, M. G. and Parr, R. Least-squares policy iteration. *The Journal of Machine Learning Research*, 4: 1107–1149, 2003.
- Le, H., Voloshin, C., and Yue, Y. Batch policy learning under constraints. In *International Conference on Machine Learning*, pp. 3703–3712, 2019.
- Liu, Q., Li, L., Tang, Z., and Zhou, D. Breaking the curse of horizon: Infinite-horizon off-policy estimation. Advances in Neural Information Processing Systems, 31: 5356–5366, 2018.
- Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. Off-policy policy gradient with state distribution correction. arXiv preprint arXiv:1904.08473, 2019.
- Nachum, O. and Dai, B. Reinforcement learning via fenchelrockafellar duality. arXiv preprint arXiv:2001.01866, 2020.
- Nachum, O., Chow, Y., Dai, B., and Li, L. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. arXiv preprint arXiv:1906.04733, 2019a.
- Nachum, O., Dai, B., Kostrikov, I., Chow, Y., Li, L., and Schuurmans, D. Algaedice: Policy gradient from arbitrary experience. arXiv preprint arXiv:1912.02074, 2019b.
- Paine, T. L., Paduraru, C., Michi, A., Gulcehre, C., Zolna, K., Novikov, A., Wang, Z., and de Freitas, N. Hyperparameter selection for offline reinforcement learning. *arXiv preprint arXiv:2007.09055*, 2020.
- Perdomo, J. C., Krishnamurthy, A., Bartlett, P., and Kakade, S. A sharp characterization of linear estimators for offline policy evaluation. arXiv preprint arXiv:2203.04236, 2022.

- Precup, D., Sutton, R. S., and Singh, S. P. Eligibility traces for off-policy policy evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pp. 759–766, 2000.
- Thomas, P. and Brunskill, E. Data-Efficient Off-Policy Policy Evaluation for Reinforcement Learning. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- Thomas, P., Theocharous, G., and Ghavamzadeh, M. High Confidence Off-Policy Evaluation. In *Proceedings of the* 29th AAAI Conference on Artificial Intelligence, 2015.
- Uehara, M., Huang, J., and Jiang, N. Minimax Weight and Q-Function Learning for Off-Policy Evaluation. In Proceedings of the 37th International Conference on Machine Learning, pp. 1023–1032, 2020.
- Uehara, M., Imaizumi, M., Jiang, N., Kallus, N., Sun, W., and Xie, T. Finite sample analysis of minimax offline reinforcement learning: Completeness, fast rates and firstorder efficiency. arXiv preprint arXiv:2102.02981, 2021.
- Voloshin, C., Le, H. M., Jiang, N., and Yue, Y. Empirical study of off-policy policy evaluation for reinforcement learning. arXiv preprint arXiv:1911.06854, 2019.
- Xie, T., Ma, Y., and Wang, Y.-X. Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. *arXiv preprint arXiv:1906.03393*, 2019.
- Xie, T., Cheng, C.-A., Jiang, N., Mineiro, P., and Agarwal, A. Bellman-consistent pessimism for offline reinforcement learning. arXiv preprint arXiv:2106.06926, 2021.
- Yang, M., Nachum, O., Dai, B., Li, L., and Schuurmans, D. Off-policy evaluation via the regularized lagrangian. *Advances in Neural Information Processing Systems*, 33: 6551–6561, 2020.
- Zhan, W., Huang, B., Huang, A., Jiang, N., and Lee, J. D. Offline reinforcement learning with realizability and singlepolicy concentrability. *arXiv preprint arXiv:2202.04634*, 2022.
- Zhang, J., Hong, M., Wang, M., and Zhang, S. Generalization bounds for stochastic saddle point problems. In *International Conference on Artificial Intelligence and Statistics*, pp. 568–576. PMLR, 2021.
- Zhang, S. and Jiang, N. Towards hyperparameter-free policy selection for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.

# A. Further Discussions

**Faster rate** One weakness of our result is the  $O(n^{-1/4})$  slow rate of estimation. While  $O(n^{-1/2})$  generalization error bounds for related stochastic saddle point exist (Zhang et al., 2021), they only apply to strongly-convex-strongly-concave problems, whereas our problem is strongly-convex-non-strongly-concave  $(L_f^q)$  is affine in w and  $L_f^w$  is affine in q), making the result not directly applicable. One immediate idea is to introduce dual regularization to make our objectives also strongly concave in the discriminator. However, while primal regularization does not change the feasible space and guarantees that the learned function will be  $q^{\pi}$  (or  $w^{\pi}$ , respectively), dual regularization *does* change the optimal solution, introducing a bias. This leads to a trade-off between the improvement in error bounds due to strong concavity and the additional bias, and our preliminary investigation shows that an optimal trade-off between the two sources of errors still leads to an  $O(n^{-1/4})$  rate. Therefore, improving the rate (if it is possible at all) will require novel technical tools for the generalization analyses of strongly-convex-non-strongly-concave stochastic saddle point problems, which will be an interesting future direction.

On a related note, while the rate for estimating  $q^{\pi}$  and  $w^{\pi}$  is only  $O(n^{-1/4})$ , we can combine them in a doubly robust form to get  $O(n^{-1/2})$  rate for return estimation by careful choices of the regularizing distributions  $\nu$  and  $\eta$ ; see Appendix F for details.

**Comparison to off-policy learning** As mentioned earlier, our results are enabled by technical tools adapted from Zhan et al. (2022), whose work focuses on off-policy policy learning and learns  $w^{\pi}$  for a near-optimal  $\pi$  that is accurate under  $d^D$  as an intermediate step. While most of our surprising observations are in the value-function learning scenario (Section 4), comparing our guarantee for learning  $w^{\pi}$  (Section 5) to that of Zhan et al. (2022) still yields interesting observations about the difference between off-policy evaluation and learning. Most notably, we do not need to control the strength of regularization in (3), since the feasible space is a singleton and there is no objective before we introduce  $\mathbb{E}_{\nu}[f(q)]$ . In contrast, the feasible space is not a singleton in Zhan et al. (2022) (it is the space of all possible occupancies) and there is already a return optimization objective, so Zhan et al. (2022) need to carefully control the strength of their regularization. As a consequence, Zhan et al. (2022) obtain  $O(n^{-1/6})$  rate, showing how off-policy learning is potentially more difficult than off-policy function estimation. Another interesting difference is related to our exact characterization of  $w_f^*$  and  $q_f^*$ : Zhan et al. (2022) do not have a closed-form expression for their optimal dual solution. Such a lack of direct characterization leads to requiring additional assumptions to guarantee the boundedness of such variables (see their Assumptions 11 and 12), which is not a problem in our setting. Finally, our analyses lead to novel algorithmic ideas such as using state-action-dependent regularizers and incorporating approximate models in the regularizers, which are potentially also useful for policy learning.

#### **B.** Proofs for Section 4

#### **B.1. Proof of Theorem 2**

From Assumption 1 and Lemma 7, we know that the regularization function  $\mathbb{E}_{\nu}[f_{s,a}(q(s,a))]$  is an *M*-strongly convex function in *q* on the  $\|\cdot\|_{2,\nu}$  norm. Now consider  $L_f^q(q, w_f^*)$ , the Lagrangian function (4) at the optimal discriminator  $w_f^*$ . Since  $L_f^q(q, w_f^*)$  is composed of the regularization function plus terms that are linear in *q*,  $L_f^q(q, w_f^*)$  is also an *M*-strongly convex function in *q*.

As  $(q^{\pi}, w_f^*)$  is the saddle point solution of  $L_f^q$ , we know  $q^{\pi} = \arg \min_q L_f^q(q, w_f^*)$ . Then from the strong convexity of  $L_f^q$ ,

$$\begin{split} ||\widehat{q} - q^{\pi}||_{2,\nu} &\leq \sqrt{\frac{2\left(L_{f}^{q}(\widehat{q}, w_{f}^{*}) - L_{f}^{q}(q^{\pi}, w_{f}^{*})\right)}{M^{q}}} \\ &\leq \sqrt{\frac{4\epsilon_{stat}^{q}}{M^{q}}}, \end{split}$$
(Lemma 9)

where  $\epsilon_{stat}^{q}$  is given in Lemma 8.

We provide the helper lemmas and their proofs below:

**Lemma 7.** Suppose  $f_{s,a} : \mathbb{R} \to \mathbb{R}$  is *M*-strongly convex. Then  $\mathbb{E}_{\nu}[f_{s,a}(q(s,a))] : \mathbb{R}^{|SA|} \to \mathbb{R}$  is *M*-strongly convex on  $\|\cdot\|_{\nu}$ .

*Proof.* From the strong convexity of  $f_{s,a}$ , for any  $x, y \in \mathbb{R}$ ,

$$f_{s,a}(x) - f_{s,a}(y) \le f'_{s,a}(x)(x-y) - \frac{M}{2}(x-y)^2$$

Then for  $q, q' \in \mathbb{R}^{|SA|}$ ,

$$\begin{split} &\mathbb{E}_{\nu}[f_{s,a}(q(s,a))] - \mathbb{E}_{\nu}[f_{s,a}(q'(s,a))] \\ &\leq \mathbb{E}_{\nu}[f'_{s,a}(q(s,a))(q(s,a) - q'(s,a))] - \mathbb{E}_{\nu}[\frac{M}{2}(q(s,a) - q'(s,a))^{2}] \\ &\leq \mathbb{E}_{\nu}[f'_{s,a}(q(s,a))(q(s,a) - q'(s,a))] - \left(\min_{s,a}\frac{M}{2}\right) \mathbb{E}_{\nu}[(q(s,a) - q'(s,a))^{2}] \\ &= \langle \nabla_{q}\mathbb{E}_{\nu}[f_{s,a}(q(s,a))], q - q'\rangle - \frac{M}{2}\mathbb{E}_{\nu}[(q(s,a) - q'(s,a))^{2}] \end{split}$$

since  $\nabla_q \mathbb{E}_{\nu}[f_{s,a}(q(s,a))] = \nu \circ f_{s,a}'(q)$ , which gives our result.

**Lemma 8.** Suppose Assumption 3 holds. Then for all  $(q, w) \in \mathcal{Q} \times \mathcal{W}$ , w.p.  $\geq 1 - \delta$ ,

$$|\widehat{L}_{f}^{q}(q,w) - L_{f}^{q}(q,w)| \le \epsilon_{stat}^{q},$$

where  $\epsilon_{stat}^q = \left(C_{\mathcal{W}}^q + (1+\gamma)C_{\mathcal{W}}^q C_{\mathcal{Q}}^q\right)\sqrt{\frac{2\log\frac{2|\mathcal{W}||\mathcal{Q}|}{\delta}}{n}}.$ 

*Proof.* From the linearity of the expectation, it is clear that  $L_f^q(q, w) = \mathbb{E}[\widehat{L}_f^q]$ . Let  $l_i = w(s_i, a_i) (r(s_i, a_i) + \gamma q(s'_i, \pi) - q(s_i, a_i))$ . From Assumption 3,

$$\begin{aligned} |l_i| &\leq \|w\|_{\infty} + (1+\gamma)\|w\|_{\infty} \|q\|_{\infty} \\ &\leq C_{\mathcal{W}}^q + (1+\gamma)C_{\mathcal{W}}^q C_{\mathcal{Q}}^q \end{aligned}$$

Then using Hoeffding's inequality with union bound, for all  $q, w \in \mathcal{Q} \times \mathcal{W}$ , w.p.  $\geq 1 - \delta$ ,

$$\left|\frac{1}{n}\sum_{i=1}^{n}l_{i} - \mathbb{E}_{d^{D}}[l_{i}]\right| \leq \left(C_{\mathcal{W}}^{q} + (1+\gamma)C_{\mathcal{W}}^{q}C_{\mathcal{Q}}^{q}\right)\sqrt{\frac{2\log\frac{2|\mathcal{W}||\mathcal{Q}|}{\delta}}{n}} = \epsilon_{stat}^{q}$$

**Lemma 9.** Under Assumptions 1, 2, 3, w.p.  $\geq 1 - \delta$ ,

$$L_f^q(\widehat{q}, w_f^*) - L_f^q(q^\pi, w_f^*) \le 2\epsilon_{stat}^q.$$

where  $\epsilon_{stat}^{q}$  is given in Lemma 8.

*Proof.* We decompose the error as follows:

$$\begin{split} L_{f}^{q}(q^{\pi}, w_{f}^{*}) - L_{f}^{q}(\widehat{q}, w_{f}^{*}) &= L_{f}^{q}(q^{\pi}, w_{f}^{*}) - L_{f}^{q}(q^{\pi}, \widehat{w}(q^{\pi})) & (1) \geq 0 \\ &+ L_{f}^{q}(q^{\pi}, \widehat{w}(q^{\pi})) - \widehat{L}_{f}^{q}(q^{\pi}, \widehat{w}(q^{\pi})) & (2) \geq -\epsilon_{stat}^{q} \\ &+ \widehat{L}_{f}^{q}(q^{\pi}, \widehat{w}(q^{\pi})) - \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) & (3) \geq 0 \\ &+ \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) - \widehat{L}_{f}^{q}(\widehat{q}, w_{f}^{*}) & (4) \geq 0 \\ &+ \widehat{L}_{f}^{q}(\widehat{q}, w_{f}^{*}) - L_{f}^{q}(\widehat{q}, w_{f}^{*}) & (5) \geq -\epsilon_{stat}^{q} \end{split}$$

Combining the terms gives the result, and we provide a brief justification for each inequality below. Terms (2) and (5) follow from Lemma 8.

Term (1)  $\geq 0$  since  $(q^{\pi}, w_f^*)$  is the saddlepoint solution.

Term (3) 
$$\geq 0$$
, since  $\widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) = \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}(\widehat{q}))$ , and  $\widehat{q} = \arg \min_{q \in \mathcal{Q}} \widehat{L}_{f}^{q}(q, \widehat{w}(q))$ , and  $q^{\pi} \in \mathcal{Q}$ .  
Term (4)  $\geq 0$  because  $w_{f}^{*} \in \mathcal{W}$ .

#### B.2. Proof of Lemma 3

Since strong duality holds, the saddle point  $(q^{\pi}, w_f^*)$  satisfies the KKT conditions. Then from stationarity, for all (s, a),

$$0 = \nu(s,a)f'_{s,a}(q^{\pi}(s,a)) + \gamma \sum_{s',a'} P^{\pi}(s,a|s',a')d^{D}(s',a')w_{f}^{*}(s',a') - d^{D}(s,a)w_{f}^{*}(s,a)$$

Writing this in matrix form, letting  $f'(q^{\pi})$  be shorthand for  $[f'_{s,a}(q^{\pi}(s,a))]_{s,a} \in \mathbb{R}^{S \times A}$ ,  $w_f^*$  must satisfy the equality:

$$(I - \gamma \widetilde{P}^{\pi})(d^D \circ w_f^*) = \nu \circ f'(q^{\pi}) \implies d^D \circ w_f^* = (I - \gamma \widetilde{P}^{\pi})^{-1} \left(\nu \circ f'(q^{\pi})\right).$$

# **B.3.** Proof of Proposition 4

Rearranging the closed form of  $w_f^*$  from Lemma 3 and taking the absolute value of both sides,

$$\begin{aligned} d^{D} \circ |w_{f}^{*}| &= |(I - \gamma \widetilde{P}^{\pi})^{-1} \left(\nu \circ f'(q^{\pi})\right)| \\ &\leq \|f'(q^{\pi})\|_{\infty} |(I - \gamma \widetilde{P}^{\pi})^{-1}\nu| \\ &= \frac{1}{1 - \gamma} \|f'(q^{\pi})\|_{\infty} \cdot d_{\nu}^{\pi} \end{aligned}$$

Then dividing both sides by  $d^D$  element-wise, this implies

$$|w_{f}^{*}| \leq \frac{1}{1-\gamma} ||f'(q^{\pi})||_{\infty} \cdot (d_{\nu}^{\pi}/d^{D})$$
$$\leq \frac{1}{1-\gamma} ||f'(q^{\pi})||_{\infty} \cdot ||d_{\nu}^{\pi}/d^{D}||_{\infty}$$

As the above inequality holds for all (s, a),

$$||w_f^*||_{\infty} \le \frac{1}{1-\gamma} ||f'(q^{\pi})||_{\infty} \cdot ||d_{\nu}^{\pi}/d^D||_{\infty}.$$

### C. Proofs for Section 5

#### C.1. Derivation of Lagrangian Objective (8)

For completeness, we demonstrate how the Lagrangian objective in (8) is derived from the constrained convex program in (7). Letting  $q \in \mathbb{R}^{S \times A}$  be the dual variable, (7) can be written in Lagrangian form and rearranged as follows:

$$\begin{split} L_{f}^{w}(w,q) &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] \\ &+ \sum_{s,a} q(s,a) \left( (1-\gamma) \mu_{0}^{\pi}(s,a) + \gamma \sum_{s',a'} P^{\pi}(s,a|s',a') d^{D}(s',a') w(s',a') - d^{D}(s,a) w(s,a) \right) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] \\ &+ \sum_{s,a} q(s,a) \left( \gamma \sum_{s',a'} P^{\pi}(s,a|s',a') d^{D}(s',a') w(s',a') - d^{D}(s,a) w(s,a) \right) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] \\ &+ \sum_{s',a'} d^{D}(s',a') w(s',a') \gamma \sum_{s,a} P^{\pi}(s,a|s',a') q(s,a) - \sum_{s,a} d^{D}(s,a) w(s,a) q(s,a) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] \\ &+ \sum_{s,a} d^{D}(s,a) w(s,a) \gamma \sum_{s',a'} P^{\pi}(s',a'|s,a) q(s',a') - \sum_{s,a} d^{D}(s,a) w(s,a) q(s,a) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] \\ &+ \sum_{s,a} d^{D}(s,a) w(s,a) (\gamma \sum_{s',a'} P^{\pi}(s',a'|s,a) q(s',a') - q(s,a)) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] \\ &+ \sum_{s,a} d^{D}(s,a) w(s,a) (\gamma \sum_{s',a'} P^{\pi}(s',a'|s,a) q(s',a') - q(s,a)) \\ &= \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + (1-\gamma) \mathbb{E}_{\mu_{0}}[q(s,\pi)] + \mathbb{E}_{d^{D}}[w(s,a)(\gamma q(s,\pi) - q(s,a))] \end{split}$$

which is exactly (8).

#### C.2. Proof of Lemma 5

From the KKT stationarity conditions:

$$0 = d^{D}(s,a) \left( \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ q_{f}^{*}(s',\pi) \right] - q_{f}^{*}(s,a) \right) - \nu(s,a) f_{s,a}'(w^{\pi}(s,a))$$

or in matrix form, letting  $f'(w^{\pi})$  be shorthand for  $[f'_{s,a}(w^{\pi}(s,a))]_{s,a} \in \mathbb{R}^{S \times A}$ ,

$$\eta \circ f'(w^{\pi}) = d^D \circ (I - \gamma P^{\pi})q_f^*$$

Then  $q_f^*$  must satisfy

$$(I - \gamma P^{\pi})q_f^* = f'(w^{\pi}) \circ \eta/d^D \implies q_f^* = (I - \gamma P^{\pi})^{-1}(f'(w^{\pi}) \circ \eta/d^D)$$

#### C.3. Proof of Theorem 6

The proof is of a similar nature as the proof of Theorem 2 (Appendix B.1). From Assumption 4 and Lemma 7, we know that that  $L_f^w(w, q_f^*)$  is an *M*-strongly convex function in w on the  $|| \cdot ||_{2,\eta}$  norm. Since  $(w^{\pi}, q_f^*)$  is the saddle point solution of  $L_f^w$ , from strong convexity we know that the error of  $\hat{w}$  is bounded as

$$\begin{split} ||\widehat{w} - w^{\pi}||_{2,d^{D}} &\leq \sqrt{\frac{2\left(L_{f}^{w}(w^{\pi}, q_{f}^{*}) - L_{f}^{w}(\widehat{w}, q_{f}^{*})\right)}{M^{w}}} \\ &\leq \sqrt{\frac{4\epsilon_{stat}^{w}}{M^{w}}} \end{split}$$
(Lemma 11),

where  $\epsilon_{stat}^{w}$  is given in Lemma 10.

**Remark 5.** In Theorem 6 of the main text, there is an additional  $O(C_f^w/\sqrt{n})$  term in the statistical error  $\epsilon_{stat}^w$ , which would arise if the regularization function  $\mathbb{E}_{\eta}[f_{s,a}(w(s,a))]$  were to be estimated from samples. However, we state early on in the paper that we assume the regularizer can be calculated exactly, as sampling is a trivial extension. Correspondingly, the correct expression for the statistical error is:

$$\epsilon_{stat}^w = (1+\gamma) C_{\mathcal{W}}^w C_{\mathcal{Q}}^w \sqrt{2\log\frac{4|\mathcal{Q}||\mathcal{W}|}{\delta}/n} + (1-\gamma) C_{\mathcal{Q}}^w \sqrt{2\log\frac{4|\mathcal{Q}|}{\delta}/n_0},$$

and, to remain consistent with the rest of the paper, we provide the proof and lemma for this  $\epsilon_{stat}^{w}$  below.

**Lemma 10.** Suppose Assumption 6 holds. Then for all  $(w, q) \in W \times Q$ , w.p.  $\geq 1 - \delta$ ,

$$|\widehat{L}_f^w(w,q) - L_f^w(w,q)| \le \epsilon_{stat}^w,$$

where  $\epsilon_{stat}^w = (1+\gamma)C_{\mathcal{W}}^w C_{\mathcal{Q}}^w \sqrt{\frac{2\log\frac{4|\mathcal{Q}||\mathcal{W}|}{\delta}}{n}} + (1-\gamma)C_{\mathcal{Q}}^w \sqrt{\frac{2\log\frac{4|\mathcal{Q}|}{\delta}}{n_0}}.$ 

*Proof.* Let  $l_i = w(s_i, a_i)(\gamma q(s'_i, \pi) - q(s_i, a_i))$ . Using Assumption 6,

$$\begin{aligned} |l_i| &\leq (1+\gamma)||w||_{\infty} ||q||_{\infty} \\ &\leq (1+\gamma)C_{\mathcal{W}}^w C_{\mathcal{Q}}^w \end{aligned}$$

Then using Hoeffding's inequality with union bound, w.p.  $\geq 1 - \delta/2$  we have that for all  $w, q \in W \times Q$ ,

$$\left|\frac{1}{n}\sum_{i=1}^{n}l_{i} - \mathbb{E}_{d^{D}}[l_{i}]\right| \leq (1+\gamma)C_{\mathcal{W}}^{w}C_{\mathcal{Q}}^{w}\sqrt{\frac{2\log\frac{4|\mathcal{W}||\mathcal{Q}|}{\delta}}{n}}$$

Similarly, for all  $q \in \mathcal{Q}$ , w.p.  $\geq 1 - \delta/2$ ,

$$\left|\frac{1}{n_0} \sum_{i=1}^{n_0} q(s_{0,i}, \pi) - \mathbb{E}_{\mu_0}[q(s_{0,i}, \pi)]\right| \le C_{\mathcal{Q}}^w \sqrt{\frac{2\log\frac{4|\mathcal{Q}|}{\delta}}{n_0}}$$

Since  $L_f^w(w,q) = \mathbb{E}_{\eta}[f_{s,a}(w(s,a))] + \mathbb{E}_{d^D}[l_i] + \mathbb{E}_{\mu_0}[q(s_0,\pi)]$ , but the first term can be calculated exactly, taking a union bound over the above two inequalities, we have that w.p.  $\geq 1 - \delta$ ,

$$|\widehat{L}_{f}^{w}(w,q) - L_{f}^{w}(w,q)| \leq (1+\gamma)C_{\mathcal{W}}^{w}C_{\mathcal{Q}}^{w}\sqrt{\frac{2\log\frac{4|\mathcal{Q}||\mathcal{W}|}{\delta}}{n}} + (1-\gamma)C_{\mathcal{Q}}^{w}\sqrt{\frac{2\log\frac{4|\mathcal{Q}|}{\delta}}{n_{0}}}$$

**Lemma 11.** Under Assumptions 4, 5, 6, w.p.  $\geq 1 - \delta$ ,

$$L_f^w(w_f^*, q_f^*) - L_f^w(\widehat{w}, q_f^*) \le 2\epsilon_{stat}^w$$

**Proof of Lemma 11** We decompose the error as follows:

$$\begin{split} L_{f}^{w}(\widehat{w},q_{f}^{*}) - L_{f}^{w}(w^{\pi},q_{f}^{*}) &= L_{f}^{w}(\widehat{w},q_{f}^{*}) - \widehat{L}_{f}(\widehat{w},q_{f}^{*}) & (1) \geq -\epsilon_{stat}^{w} \\ &+ \widehat{L}_{f}^{w}(\widehat{w},q_{f}^{*}) - \widehat{L}_{f}^{w}(\widehat{w},\widehat{q}) & (2) \geq 0 \\ &+ \widehat{L}_{f}^{w}(\widehat{w},\widehat{q}) - \widehat{L}_{f}^{w}(w^{\pi},\widehat{q}(w^{\pi})) & (3) \geq 0 \\ &+ \widehat{L}_{f}^{w}(w^{\pi},\widehat{q}(w^{\pi})) - L_{f}^{w}(w^{\pi},\widehat{q}(w^{\pi})) & (4) \geq -\epsilon_{stat}^{w} \\ &+ L_{f}^{w}(w^{\pi},\widehat{q}(w^{\pi})) - L_{f}^{w}(w^{\pi},q_{f}^{*}) & (5) \geq 0 \end{split}$$

Combining the inequalities gives the result. We give a brief justification for each term below. Terms (1) and (4) follow from Lemma 10.

Term (2)  $\geq 0$ , since  $q_f^* \in Q$ .

Term (3)  $\geq 0$  since  $w^{\pi} \in \mathcal{W}$  and  $\widehat{w} = \arg \max_{w \in \mathcal{W}} \widehat{L}_{f}^{w}(w, \widehat{q}(w)).$ 

Term (5)  $\geq 0$  since  $(w^{\pi}, q_f^*)$  is a saddle point solution.

# **D.** Additional Details of the Experiments

# **D.1.** Derivation

We now derive the system of equations for our value function estimation experiments in Section 6. Letting the regularization function be  $f_{s,a}(x) = \frac{1}{2}x^2$  for all (s, a), the objective is

$$\min_{q} \max_{w} L_{f}^{q}(q,w) = \frac{1}{2} \mathbb{E}_{\nu}[q^{2}(s,a)] + \mathbb{E}_{d^{D}}\left[w(s,a)\left(r(s,a) + \gamma q(s',\pi) - q(s,a)\right)\right],$$
(10)

Letting  $\mathbb{E}_n$  denote the empirical average over  $\mathcal{D}$  for clarity, with empirical samples and the linear classes  $\mathcal{Q}, \mathcal{W}$ , the objective becomes:

$$\min_{q \in \mathcal{Q}} \max_{w \in \mathcal{W}} \widehat{L}_{f}^{q}(q, w) = \frac{1}{2} \mathbb{E}_{\nu} [\alpha^{\top} \phi(s, a) \phi(s, a)^{\top} \alpha] + \beta^{\top} \Big( \mathbb{E}_{n} \left[ \phi(s, a) r(s, a) \right] \\ + \mathbb{E}_{n} \left[ \gamma \phi(s, a) \phi(s', \pi)^{\top} - \phi(s, a) \phi(s, a)^{\top} \right] \alpha \Big)$$

Since  $\beta \in \mathbb{R}^d$ ,  $\max_{w \in \mathcal{W}} \widehat{L}_f^q(q, w) = +\infty$  for any q, unless  $\alpha$  sets the second term to 0. This is satisfied by  $\alpha$  such that

$$\mathbb{E}_n\left[\phi(s,a)\phi(s,a)^\top - \gamma\phi(s,a)\phi(s',\pi)^\top\right]\alpha = \mathbb{E}_n\left[\phi(s,a)r(s,a)\right].$$

However, there may in general be infinite feasible  $\alpha$  depending on the linear features and samples. For our specific linear parameterization of Q, W, the constraints form an underdetermined  $d \times k$  system of equations, which has infinite solutions.

This is where the regularization term  $\mathbb{E}_{\nu}[\alpha^{\top}\phi(s,a)\phi(s,a)^{\top}\alpha]$  comes into play. For any regularizing distribution  $\nu$ , our method will output a solution that minimizes this term, i.e. that minimizes the norm of  $q = \Phi^{\top}\alpha$  on  $\nu$ . If  $\nu = 0$ , for example, the algorithm will output any feasible point; if  $\nu = 1/|\mathcal{SA}|$ , the algorithm will output q with smallest L2 norm.

**Connection to LSTDQ** When using the same linear class for W and Q, the solution to the constraints in (3) (i.e., ignoring the regularization objective)—if the solution is unique given matrix invertibility—coincides with LSTDQ (Uehara et al., 2020). As mentioned in Section 2, LSTDQ enjoys function-estimation guarantees under matrix invertibility. In fact, we believe it is possible to extend the analysis even when Q and W use different features of dimensions d and k, respectively; as long as  $k \ge d$  and the matrix in (3) has full row-rank<sup>8</sup> (i.e., *overdetermined*), similar guarantees for LSTDQ should still hold, though we are not aware of an explicit documentation of this fact. In contrast, our setup is more challenging as we are in the regime of k < d, and the constraints in (3) is *underdetermined*, nullifying the guarantees of LSTDQ. In such cases, the use of regularization is important for guaranteeing function estimation, as also shown in our experiments.

#### **D.2. Experimental Setup**

**Feature Design** In total, the tabular environment has 400 state-action values, and we design  $\Phi$  to aggregate states that correspond to unique entries (within 3 decimal places) of  $q^{\pi}$ . In Figure 1,  $\tilde{\Phi}$  is composed of the set of features given by

$$\{(I - \gamma \widetilde{P}^{\pi})^{-1}(\nu \circ q^{\pi})/d^{D}, (I - \gamma \widetilde{P}^{\pi})^{-1}(\nu \circ q^{\pi})\}_{\nu \in \mathcal{V}}$$

The first of these two entries is the closed-form solution of  $w_f^*$  given in Lemma 3, and satisfies the realizability requirements of all methods; the second is included for optimization stability.

In Figure 2, we use a model with constant value equal to the average value of  $q^{\pi}$  on the support of p, i.e.  $\overline{q} = \frac{1}{|S\mathcal{A}| \sum_{s,a} q^{\pi}(s,a) \cdot \mathbb{1}_{\{p>0\}}}$ . To maintain realizability when the model is included in the regularization function,  $\widetilde{\Phi}$  is composed of the set

$$\left\{ (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ q^{\pi}), \ (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ q^{\pi} \circ \mathbb{I}(\widetilde{q} > 0)), \ (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ \mathbb{I}(\widetilde{q} > 0)) \right\}_{\nu \in \mathcal{V}}$$

The reason why this preserves realizability is as follows. When  $\nu$  is the regularization distribution, and the input model is  $\tilde{q} = (mq^{\pi} + (1-m)\bar{q}) \circ \mathbb{1}(p > 0))$  for some constant  $\bar{q}$ , the closed-form solution  $w_f^*$  can be expanded as

$$w_f^* = (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ (q^{\pi} - \widetilde{q}))$$
  
=  $(I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ q^{\pi}) - m \cdot (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ \mathbb{1}(p > 0) \circ q^{\pi})$   
 $- (1 - m)\overline{q} \cdot (I - \gamma \widetilde{P}^{\pi})^{-1} (\nu \circ \mathbb{1}(p > 0)),$ 

which implies  $w_f^*$  can be expressed as a linear combination of the three previously defined features.

Solver We solve the linear system using CVXPY with optimizer SCS (Diamond & Boyd, 2016; Agrawal et al., 2018).

**Environment** The Gridwalk is a 10x10 environment with 4 actions corresponding to cardinal directions. The objective is to reach the goal state (lower right corner). In each state, the agent receives a reward inversely proportional to its distance from a goal state. Each trajectory terminates after 100 steps. The initial states are randomly distributed over the upper half of the grid.

The target policy is defined to be a deterministic optimal policy that always moves towards the goal by first going right, and then down. To create a strong shift, the behavioral policy is designed to largely explore only the bottom left portion of the grid, providing poor coverage over the target policy and starting states. Specifically, letting the following probabilities refer to distributions over actions [RIGHT, DOWN, LEFT, UP], the target policy  $\pi$  has distribution [1, 0, 0, 0] over actions until it hits the right wall, then [0, 1, 0, 0]. The behavior policy takes [0.1, 0.4, 0.5, 0] until it hits the right wall, then takes [0, 0.5, 0.5, 0].

# E. Approximation and Optimization Error

The main results of this paper (Theorems 2, 6) utilize assumptions on realizability (Assumption 2, 5), as well as (implicit) assumptions of perfect optimization. In this section, we analyze how approximation errors, i.e. when the saddle point solution is not contained in  $Q \times W$ , and optimization errors affect our error bounds. Due to the similarity in proofs between value function and weight learning, we provide them only for value function learning; analogous methods can be used to derive similar results for weight learning.

<sup>&</sup>lt;sup>8</sup> In the finite-sample regime, one needs to lower-bound the smallest singular value of such matrices instead of imposing full-rankness (Perdomo et al., 2022).

#### **E.1. Finite-sample Guarantees**

First, we relax the realizability requirements of Assumption 2. Define the approximation errors:

$$\begin{aligned} \epsilon_{approx,q} &= \min_{q \in \mathcal{Q}} \max_{w \in \mathcal{W}} |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a))] + \mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))]| \\ \epsilon_{approx,w} &= \min_{w \in \mathcal{W}} \max_{q \in \mathcal{Q}} |\mathbb{E}_{d^{D}}[(w(s,a) - w_{f}^{*}(s,a))(\mathcal{T}^{\pi}q(s,a) - q(s,a))]| \\ \epsilon_{approx} &:= \epsilon_{approx,q} + \epsilon_{approx,w}. \end{aligned}$$

 $\epsilon_{approx,q}$  is composed of the worst-case weighted combination of Bellman errors of the best candidate  $q \in Q$ , as well as the difference between the regularization function at q and  $q^{\pi}$ . The error  $\epsilon_{approx,w}$  measures the distance between the best candidate  $w \in W$  and the saddle point solution  $w_f^*$  by projecting the difference onto the worst-case Bellman error  $\mathcal{T}^{\pi}q - q$ .

**Remark 6.** To increase intuition of  $\epsilon_{approx,q}$ , we can relax the difference in regularization terms as  $\mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))] \leq C_{f'}^{q}||q^{\pi} - q||_{2,\nu}$ , which is also the norm upon which the  $\hat{q}$  estimation guarantee is given (Theorem 2). Reflecting the nature of the value function estimation task, this states that, even if there is a candidate  $q \in \mathcal{Q}$  with low Bellman error (e.g. if data is sparse),  $\epsilon_{approx,q}$  will still be large if q is far from  $q^{\pi}$  on the desired distribution  $\nu$ .

Next, we can also relax the (implicit) assumptions that the estimates  $(\hat{q}, \hat{w})$  are the true optima of (5), e.g.  $\hat{q} = \arg \min_{q \in \mathcal{Q}} \hat{L}_{f}^{q}(q, \hat{w}(q))$  and  $\hat{w} = \arg \max_{w \in \mathcal{W}} \hat{L}_{f}^{q}(\hat{q}, w)$ . Letting  $\hat{w}(q) = \arg \max_{w \in \mathcal{W}} \hat{L}(w, q)$ , define the following optimization errors:

$$\begin{split} \epsilon_{opt,w} &\geq \hat{L}_{f}^{q}(\hat{q}, \hat{w}(\hat{q})) - \hat{L}_{f}^{q}(\hat{q}, \hat{w}) \\ \epsilon_{opt,q} &\geq \hat{L}_{f}^{q}(\hat{q}, \hat{w}(\hat{q})) - \min_{q \in \mathcal{Q}} \hat{L}_{f}^{q}(q, \hat{w}(q)) \\ \epsilon_{opt} &:= \epsilon_{opt,q} + \epsilon_{opt,w}. \end{split}$$

 $\epsilon_{opt,w}$  states that the estimate  $\hat{w}$  should not be too far from the best discriminator in  $\mathcal{W}$  for  $\hat{q}$ , while  $\epsilon_{opt,q}$  states that the estimate  $\hat{q}$  should not be too far from the minimax solution.

Using the above definitions, we provide the following generalization of Theorem 2, which accounts for approximation and optimization errors.

**Theorem 12.** Under Assumptions 1 and 3, with probability at least  $1 - \delta$ ,

$$||\widehat{q} - q^{\pi}||_{2,\nu} \le \sqrt{\frac{4\epsilon_{stat}^{q} + 2\epsilon_{approx} + 2\epsilon_{opt}}{M^{q}}},$$

where  $\epsilon_{stat}^{q}$  is given in Theorem 2.

#### E.2. Proof of Theorem 12

The proof takes the same overall steps as the proof of Theorem 2 (Appendix B.1), but relies on Lemma 13 to incorporate the approximation and optimization errors:

$$\begin{split} ||\widehat{q} - q^{\pi}||_{2,\nu} &\leq \sqrt{\frac{2\left(L_{f}^{q}(\widehat{q}, w_{f}^{*}) - L_{f}^{q}(q^{\pi}, w_{f}^{*})\right)}{M^{q}}} \\ &\leq \sqrt{\frac{4\epsilon_{stat}^{q} + 2\epsilon_{approx,q} + 2\epsilon_{approx,w} + 2\epsilon_{opt,q} + 2\epsilon_{opt,w}}{M^{q}}}. \end{split}$$
(Lemma 13)

Below, we state and prove the helper lemma, which bounds the difference between the Lagrangian objective (4) at the saddle point  $(q^{\pi}, w_f^*)$  and the point  $(\hat{q}, w_f^*)$ :

**Lemma 13.** Under Assumptions 1 and 3, w.p.  $\geq 1 - \delta$ ,

$$L_f^q(\widehat{q}, w_f^*) - L_f^q(q^{\pi}, w_f^*) \le 2\epsilon_{stat}^q + \epsilon_{approx,q} + \epsilon_{approx,w} + \epsilon_{opt,q} + \epsilon_{opt,w}.$$

*Proof.* With some abuse of notation (as  $\tilde{q}, \tilde{w}$  previously referred to models used with the regularizer), for brevity in this section, let  $\tilde{q}$  be the minimizer of  $\epsilon_{approx,q}$  and  $\tilde{w}$  be the minimizer of  $\epsilon_{approx,w}$ . That is,

$$\begin{split} \widetilde{q} &= \underset{q \in \mathcal{Q}}{\arg\min} \max_{w \in \mathcal{W}} |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a))] + \mathbb{E}_{\nu}[f(q(s,a)) - f(q^{\pi}(s,a))]|\\ \widetilde{w} &= \underset{w \in \mathcal{W}}{\arg\min} \max_{q \in \mathcal{Q}} |\mathbb{E}_{d^{D}}[(w(s,a) - w_{f}^{*}(s,a))(\mathcal{T}^{\pi}q(s,a) - q(s,a))]|. \end{split}$$

Decompose the error as follows:

$$\begin{split} L_{f}^{q}(q^{\pi}, w_{f}^{*}) - L_{f}^{q}(\widehat{q}, w_{f}^{*}) &= L_{f}^{q}(q^{\pi}, w_{f}^{*}) - L_{f}^{q}(q^{\pi}, \widehat{w}(\widetilde{q})) & (1) \geq 0 \\ &+ L_{f}^{q}(q^{\pi}, \widehat{w}(\widetilde{q})) - L_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) & (2) \geq -\epsilon_{approx,q} \\ &+ L_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) - \widehat{L}_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) & (3) \geq -\epsilon_{stat} \\ &+ \widehat{L}_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) - \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) & (4) \geq -\epsilon_{opt,q} \\ &+ \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) - \widehat{L}_{f}^{q}(\widehat{q}, \widetilde{w}) & (5) \geq -\epsilon_{opt,w} \\ &+ \widehat{L}_{f}^{q}(\widehat{q}, \widetilde{w}) - L_{f}^{q}(\widehat{q}, \widetilde{w}) & (6) \geq -\epsilon_{stat} \\ &+ L_{f}^{q}(\widehat{q}, \widetilde{w}) - L_{f}^{q}(\widehat{q}, w_{f}^{*}) & (7) \geq -\epsilon_{approx,w} \end{split}$$

First, (1) holds because  $(q^{\pi}, w_f^*)$  is the saddle point solution of  $L_f^q$  over all  $q, w \in \mathbb{R} \times \mathbb{R}$ . The statistical errors in (3) and (6) follow from Lemma 8.

Next, we justify the optimization errors. For (4),

$$\widehat{L}_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) - \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}) \geq \widehat{L}_{f}^{q}(\widetilde{q}, \widehat{w}(\widetilde{q})) - \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}(\widehat{q})) \geq \min_{q \in \mathcal{Q}} \widehat{L}_{f}^{q}(q, \widehat{w}(q)) - \widehat{L}_{f}^{q}(\widehat{q}, \widehat{w}(\widehat{q})) \geq -\epsilon_{opt,q}.$$

For (5),

$$\widehat{L}_{f}^{q}(\widehat{q},\widehat{w}) - \widehat{L}_{f}^{q}(\widehat{q},\widetilde{w}) \geq \widehat{L}_{f}^{q}(\widehat{q},\widehat{w}) - \max_{w \in \mathcal{W}} \widehat{L}_{f}^{q}(\widehat{q},w) \geq -\epsilon_{opt,w}$$

Finally, we justify the approximation errors, starting with (2). Note that for any  $q, w \in \mathcal{Q} \times \mathcal{W}$ ,

$$\begin{split} |L_{f}^{q}(q^{\pi},w) - L_{f}^{q}(q,w)| \\ &= |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a) - \mathcal{T}^{\pi}q^{\pi}(s,a) + q^{\pi}(s,a))] \\ &+ \mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))]| \\ &= |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a))] + \mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))]| \\ &\leq \max_{w \in \mathcal{W}} |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a))] + \mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))]|. \end{split}$$

Then since  $\tilde{q}$  was chosen to minimize the above expression,

$$\begin{split} L_{f}^{q}(q^{\pi},\widehat{w}(\widetilde{q})) &- L_{f}^{q}(\widetilde{q},\widehat{w}(\widetilde{q}))\\ \geq &- \max_{w\in\mathcal{W}} |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}\widetilde{q}(s,a) - \widetilde{q}(s,a))] + \mathbb{E}_{\nu}[f_{s,a}(\widetilde{q}(s,a)) - f_{s,a}(q^{\pi}(s,a))]|\\ &= &- \min_{q\in\mathcal{Q}} \max_{w\in\mathcal{W}} |\mathbb{E}_{d^{D}}[w(s,a)(\mathcal{T}^{\pi}q(s,a) - q(s,a))] + \mathbb{E}_{\nu}[f_{s,a}(q(s,a)) - f_{s,a}(q^{\pi}(s,a))]|\\ &= &-\epsilon_{approx,q}. \end{split}$$

Next we justify (8). For any  $w \in W$  and  $q \in Q$ ,

$$\begin{aligned} |L_{f}^{q}(q,w) - L_{f}^{q}(q,w_{f}^{*})| &= |\mathbb{E}_{d^{D}}[(w(s,a) - w_{f}^{*}(s,a))(\mathcal{T}^{\pi}q(s,a) - q(s,a))]| \\ &\leq \max_{q \in \mathcal{Q}} |\mathbb{E}_{d^{D}}[(w(s,a) - w_{f}^{*}(s,a))(\mathcal{T}^{\pi}q(s,a) - q(s,a))]|. \end{aligned}$$

Then since  $\tilde{w}$  was chosen to minimize the RHS of the above inequality,

$$\begin{split} L_f^q(\widehat{q}, \widetilde{w}) - L_f^q(\widehat{q}, w_f^*) &\geq -\max_{q \in \mathcal{Q}} |\mathbb{E}_{d^D}[(\widetilde{w}(s, a) - w_f^*(s, a))(\mathcal{T}^{\pi}q(s, a) - q(s, a))]| \\ &= -\min_{w \in \mathcal{W}} \max_{q \in \mathcal{Q}} |\mathbb{E}_{d^D}[(w(s, a) - w_f^*(s, a))(\mathcal{T}^{\pi}q(s, a) - q(s, a))]| \\ &= -\epsilon_{approx, w}. \end{split}$$

Combining these inequalities gives the lemma statement.

r \_ /

#### F. Off-Policy Return Estimation

. \_\_\_\_

Section 4 demonstrates how q-value estimates  $\hat{q}$  can be obtained, and Section 5 demonstrates how weight estimates  $\hat{w}$  can be obtained. The estimates  $\hat{q}$  and/or  $\hat{w}$  can additionally be used for downstream off-policy evaluation (OPE) of the policy's value  $J(\pi)$ , which can be equivalently defined in the following three ways:

$$J(\pi) = (1 - \gamma) \mathbb{E}_{s_0 \sim \mu_0} [q^{\pi}(s_0, \pi)]$$

$$J(\pi) = \mathbb{E}_{(s,a) \sim d^D, r \sim R(\cdot|s,a)} [w^{\pi}(s,a) \cdot r]$$

$$J(\pi) = (1 - \gamma) \mathbb{E}_{s_0 \sim \mu_0} [q^{\pi}(s_0, \pi)]$$

$$+ \mathbb{E}_{(s,a) \sim d^D, r \sim R(\cdot|s,a), s' \sim P(\cdot|s,a)} [w^{\pi}(s,a)(r + q^{\pi}(s', \pi) - q^{\pi}(s,a))]$$
("value function-based")
("weight-based")
("doubly robust")

With finite samples and estimates  $\hat{q}$  and  $\hat{w}$  approximating  $q^{\pi}$  and  $w^{\pi}$ , respectively, their corresponding off-policy estimators are:

$$\begin{aligned} \widehat{J}^{q}(\pi) &= (1-\gamma)\frac{1}{n_{0}}\sum_{i=1}^{n_{0}}\widehat{q}(s_{0,i},\pi) \\ \widehat{J}^{w}(\pi) &= \frac{1}{n}\sum_{i=1}^{n}\widehat{w}(s_{i},a_{i})r_{i} \\ \widehat{J}^{dr}(\pi) &= (1-\gamma)\frac{1}{n_{0}}\sum_{j=1}^{n_{0}}\widehat{q}(s_{0,j},\pi) + \frac{1}{n}\sum_{i=1}^{n}\widehat{w}(s_{i},a_{i})\left(r_{i}+\widehat{q}(s_{i}',\pi)-\widehat{q}(s_{i},a_{i})\right) \end{aligned}$$

While the OPE estimator  $\hat{J}^{dr}(\pi)$  utilizes both the weights and value functions,  $\hat{J}^w(\pi)$  and  $\hat{J}^q(\pi)$  utilize only one or the other. As a result, when  $\hat{q}$  and  $\hat{w}$  are estimated as in Sections 4 and 5, respectively,  $\hat{J}^w(\pi)$  and  $\hat{J}^q(\pi)$  both inherit their  $O(n^{-1/4})$  sample complexities:

**Corollary 14.** Suppose Assumptions 1, 2, and 3 hold, and let  $(\hat{q}, \_) = \arg \min_{q \in \mathcal{Q}} \arg \max_{w \in \mathcal{W}} \hat{L}_{f}^{q}(q, w)$ . Then with probability  $\geq 1 - 2\delta$ ,

$$|\widehat{J}^q(\pi) - J(\pi)| \le \epsilon_{eval}^q + \sqrt{\mathcal{C}_{\mu_0^\pi/\nu}} \cdot \epsilon_{est}^q,$$

where  $\epsilon_{eval}^q = (1 - \gamma) C_{\mathcal{Q}}^q \sqrt{2 \log \frac{2|\mathcal{Q}|}{\delta} / n_0}$ ,  $\mathcal{C}_{\mu_0^\pi/\nu} = ||\mu_0^\pi/\nu||_{\infty}$ , and  $\epsilon_{stat}^q$  is as in Theorem 2.

**Corollary 15.** Suppose Assumption 4, 5, and 6 hold, and let  $(\widehat{w}, \_) = \arg \min_{w \in \mathcal{W}} \arg \max_{q \in \mathcal{Q}} \widehat{L}_{f}^{w}(q, w)$ . Then with probability  $\ge 1 - 2\delta$ ,

$$|\widehat{J}^w(\pi) - J(\pi)| \le \epsilon^w_{eval} + \sqrt{\mathcal{C}_{d^D/\eta}} \cdot \epsilon^w_{est},$$

where  $\epsilon_{eval}^w = C_W^w \sqrt{\frac{2\log \frac{2|W|}{\delta}}{n}}$ ,  $C_{d^D/\eta} = \|d^D/\eta\|_{\infty}$ , and  $\epsilon_{est}^w$  is as in Theorem 6.

However, when  $\hat{q}$  and  $\hat{w}$  are used together in the doubly robust estimator  $\hat{J}^{dr}$ , their estimation error becomes multiplicative, and  $\hat{J}^{dr}(\pi)$  can achieve the  $O(n^{-\frac{1}{2}})$  fast rate of convergence. In Theorem 16 below, we present two versions this guarantee. The first requires no additional assumptions beyond  $d^D > 0$ , which we already make (see footnote 5), but involves the largest singular value of  $I - \gamma P^{\pi}$ , which may be difficult to characterize. The second utilizes an additional assumption, and replaces the singular value with an occupancy ratio, stated below. The assumption requires that all next states s' are also present as states s in transitions of  $d^D$  (a condition which may reasonably hold in practice), and is also made by (Uehara et al., 2021).

Assumption 7 (Next State Coverage). Let  $d^D(s) = \sum_a d^D(s, a)$  be the marginal distribution of states s in  $d^D$ , and  $d^D_{s'}(s) := \sum_{s',a'} P(s|s',a') d^D(s',a')$  be the marginal distribution of next states s'. Suppose

$$\mathcal{C}_{s'/s} := ||d_{s'}^D(\cdot)/d^D(\cdot)||_{\infty} < \infty$$

**Theorem 16.** Suppose Assumption 1, 2, 3, 4, 5, and 6 hold. Let  $\hat{w}$  and  $\hat{q}$  be estimated from:

$$\begin{split} & (\widehat{q}, \ \_) = \mathop{\arg\min}_{q \in \mathcal{Q}} \mathop{\arg\max}_{w \in \mathcal{W}} \widehat{L}_{f}^{q}(q, w) \\ & (\widehat{w}, \ \_) = \mathop{\arg\min}_{w \in \mathcal{W}} \mathop{\arg\max}_{q \in \mathcal{Q}} \widehat{L}_{f}^{w}(q, w). \end{split}$$

Then with probability  $\geq 1 - 3\delta$ ,

$$|\widehat{J}^{dr}(\pi) - J(\pi)| \le \epsilon_{eval}^{dr} + \sigma_{max}(I - \gamma P^{\pi}) \cdot \sqrt{\mathcal{C}_{d^D/\eta}\mathcal{C}_{d^D/\nu}} \cdot \epsilon_{est}^w \cdot \epsilon_{est}^q,$$

If Assumption 7 additionally holds, with probability  $\geq 1 - 3\delta$ ,

$$\widehat{J}^{dr}(\pi) - J(\pi)| \le \epsilon_{eval}^{dr} + \left(1 + \gamma \sqrt{\mathcal{C}_{s'/s} \mathcal{C}_{\pi/\pi^D}}\right) \cdot \sqrt{\mathcal{C}_{d^D/\eta} \mathcal{C}_{d^D/\nu}} \cdot \epsilon_{est}^w \cdot \epsilon_{est}^q$$

where  $\epsilon_{eval}^{dr} = (1-\gamma)C_{Q}^{q}\sqrt{\frac{2\log\frac{2|Q|}{\delta}}{n_0}} + C_{W}^{w}(1+(1+\gamma)C_{Q}^{q})\sqrt{\frac{2\log\frac{2|W||Q|}{\delta}}{n}}, \sigma_{max}$  denotes the largest singular value, and  $\epsilon_{est}^{q}$  and  $\epsilon_{est}^{w}$  are as in Theorems 2 and 6.

As the evaluation error  $\epsilon_{eval}^{dr}$  in Theorem 16 is  $O(n^{-1/2})$ , the sample complexity of doubly robust estimation is rate-limited by  $\epsilon_{est}^w \cdot \epsilon_{est}^q$ , the product of weight and value function estimation errors. If both functions can be estimated at an  $O(n^{-1/4})$ rate, as is true of our method, then  $\hat{J}^{dr}(\pi)$  attains the overall  $O(n^{-1/2})$  fast rate. Finally, while Theorem 16 assumes for simplicity that the same Q, W classes are used in both of its optimization problems, it can easily be extended to the case where different pairs of function classes are used as long as the required assumptions hold.

**Remark 7** (Comparison to Related Work). (Yang et al., 2020) conduct experiments comparing off-policy evaluation using  $\hat{J}^q(\pi)$ ,  $\hat{J}^w(\pi)$ ,  $\hat{J}^{dr}(\pi)$ , and generally observe that  $\hat{J}^{dr}(\pi)$  has higher variance and worse performance than either  $\hat{J}^q(\pi)$  or  $\hat{J}^w(\pi)$ . Though at first glance this may appear to contradict Theorem 16, that is actually not the case; in fact, our theoretical analysis provides insight into why (Yang et al., 2020) may observe such a phenomenon. In contrast to Theorem 16, when using  $\hat{J}^{dr}(\pi)$  (Yang et al., 2020) utilize saddle point predictions  $(\hat{q}, \hat{w})$  from either *only* value function learning or *only* weight learning, e.g.  $(\hat{q}, \hat{w}) = \arg\min_{q \in Q} \arg\max_{w \in W} \hat{L}_f^q(q, w)$  that approximates  $(q^{\pi}, w_f^*)$ . Continuing with this example (and the same applies to weight learning), it is clear from our analysis that  $\hat{w}$  estimated in such a manner may not approximate  $w^{\pi}$  at all, leading to increased estimation error of  $\hat{J}^{dr}(\pi)$  over  $\hat{J}^q(\pi)$ . First, the closed-form solution we have derived for  $w_f^*$  in (Lemma 3) shows that  $w_f^*$  may have a significantly different magnitude from  $w^{\pi}$ . Second, even if  $\nu$  and f were chosen such that  $w_f^* \approx w^{\pi}$ , as per the reasons stated in Section 4.1, we are not even guaranteed to output  $\hat{w}$  close to  $w_f^*$  since  $L_f^q$  is not regularized in w. In order to obtain the estimation benefits of doubly robust estimation, our analysis shows that  $\hat{q}$  and  $\hat{w}$  should be separately estimated from their respective optimization problems, then combined in  $\hat{J}^{dr}(\pi)$ . This is in accordance with similar results from Kallus & Uehara (2020) and Uehara et al. (2021).

#### F.1. Proof of Corollary 14

Let  $\widetilde{J}(\pi) = (1 - \gamma) \mathbb{E}_{\mu_0}[\widehat{q}(s, \pi)]$ . We decompose the error as

$$|\widehat{J}(\pi) - J(\pi)| \le |\widehat{J}(\pi) - \widetilde{J}(\pi)| + |\widetilde{J}(\pi) - J(\pi)|$$

First we bound  $|\widehat{J}(\pi) - \widetilde{J}(\pi)|$ . Using Hoeffding's with union bound, for all  $q \in \mathcal{Q}$ , w.p.  $\geq 1 - \delta$ ,

$$\left|\frac{1}{n_0} \sum_{i=1}^n q(s_{0,i}, \pi) - \mathbb{E}_{\mu_0}[q(s, \pi)]\right| \le (1 - \gamma) C_{\mathcal{Q}}^q \sqrt{\frac{2\log\frac{2|\mathcal{Q}|}{\delta}}{n_0}} := \epsilon_{eval}^q,$$

which implies  $|\widehat{J}(\pi) - \widetilde{J}(\pi)| \le \epsilon_{eval}^q$ . For the second term, let  $C_{\mu_0^\pi/\nu} = ||\mu_0^\pi/\nu||_{\infty}$ . Then w.p.  $\ge 1 - \delta$ 

$$\begin{aligned} |J(\pi) - J(\pi)| &= (1 - \gamma) |\langle \mu_0^{\pi}, \hat{q} - q^{\pi} \rangle| \\ &\leq (1 - \gamma) ||\hat{q} - q^{\pi}||_{1,\mu_0^{\pi}} \\ &\leq (1 - \gamma) ||\hat{q} - q^{\pi}||_{2,\mu_0^{\pi}} \\ &= (1 - \gamma) \sqrt{\mathcal{C}_{\mu_0^{\pi}/\nu}} ||\hat{q} - q^{\pi}||_{2,\mu_0^{\pi}} \\ &\leq (1 - \gamma) \sqrt{\mathcal{C}_{\mu_0^{\pi}/\nu}} \epsilon_{est}^q \end{aligned}$$

using Theorem 2 in the last line.

#### F.2. Proof of Corollary 15

Let  $\widetilde{J}(\pi) = \mathbb{E}_{d^D}[\widehat{w}(s, a)r(s, a)]$ . We decompose the error as

$$|\widehat{J}^w(\pi) - J(\pi)| \le |\widehat{J}^w(\pi) - \widetilde{J}(\pi)| + |\widetilde{J}(\pi) - J(\pi)|$$

For the first term, using Hoeffding's with union bound, w.p.  $\geq 1 - \delta$ , for all  $w \in \mathcal{W}$ ,

$$\left|\frac{1}{n}\sum_{i=1}^{n}w(s_i,a_i)r_i - \mathbb{E}_{d^D}[w(s,a)r(s,a)]\right| \le C_{\mathcal{W}}^w \sqrt{\frac{2\log\frac{2|\mathcal{W}|}{\delta}}{n}} := \epsilon_{eval}^w$$

which implies  $|\widehat{J}(\pi) - \widetilde{J}(\pi)| \leq \epsilon^w_{eval}.$  For the second term,

$$\begin{split} |\widehat{J}(\pi) - J(\pi)| &= |\langle \widehat{w} \cdot d^{D}, r \rangle - \langle w^{\pi} \cdot d^{D}, r \rangle| \\ &\leq ||d^{D} \cdot (\widehat{w} - w^{\pi})||_{1} ||r||_{\infty} \\ &\leq ||d^{D} \cdot (\widehat{w} - w^{\pi})||_{1} = ||\widehat{w} - w^{\pi}||_{d^{D}, 1} \\ &\leq ||\widehat{w} - w^{\pi}||_{d^{D}, 2} \\ &\leq \sqrt{\mathcal{C}_{d^{D}/\eta}} ||\widehat{w} - w^{\pi}||_{2, \eta} \\ &\leq \sqrt{\mathcal{C}_{d^{D}/\eta}} \epsilon^{w}_{est} \end{split}$$

w.p.  $\geq 1 - \delta$ , using Theorem 6 in the last line. Taking a union bound over both terms gives the stated result.

#### F.3. Proof of Theorem 16

Let  $\widetilde{J}(\pi) = (1 - \gamma) \mathbb{E}_{\mu_0^{\pi}}[\widehat{q}(s, a)] + \mathbb{E}_{d^D}[\widehat{w}(s, a)(r + \widehat{q}(s', \pi) - \widehat{q}(s, a))].$  Again we decompose the error as:  $|\widehat{J}^{dr}(\pi) - J(\pi)| \le |\widehat{J}^{dr}(\pi) - \widetilde{J}(\pi)| + |\widetilde{J}(\pi) - J(\pi)|.$ 

For the first term, since  $\mathbb{E}[\widehat{J}^{dr}(\pi)] = \widetilde{J}(\pi)$ , w.p.  $\geq 1 - \delta$  we have that  $\forall q, w \in \mathcal{Q} \times \mathcal{W}$ ,

$$|\widehat{J}^{dr}(\pi) - \widetilde{J}(\pi)| \le (1 - \gamma)C_{\mathcal{Q}}^q \sqrt{\frac{2\log\frac{2|\mathcal{Q}|}{\delta}}{n_0}} + C_{\mathcal{W}}^w (1 + (1 + \gamma)C_{\mathcal{Q}}^q) \sqrt{\frac{2\log\frac{2|\mathcal{W}||\mathcal{Q}|}{\delta}}{n}} := \epsilon_{eval}^{dr}$$

For the second term,

$$\begin{split} |\widetilde{J}(\pi) - J(\pi)| &= |(1-\gamma)\langle \widehat{q}, \mu_0^{\pi} \rangle + \langle \widehat{w} \cdot d^D, r + \gamma P^{\pi} \widehat{q} - \widehat{q} \rangle - (1-\gamma)\langle q^{\pi}, \mu_0^{\pi} \rangle| \\ &= |(1-\gamma)\langle \widehat{q}, \mu_0^{\pi} \rangle + \langle \widehat{w} \cdot d^D, r + \gamma P^{\pi} \widehat{q} - \widehat{q} \rangle - (1-\gamma)\langle q^{\pi}, \mu_0^{\pi} \rangle - \langle \widehat{w} \cdot d^D, r + \gamma P^{\pi} q^{\pi} - q^{\pi} \rangle| \\ &= |\langle \widehat{q} - q^{\pi}, (1-\gamma)\mu_0^{\pi} + (\gamma P^{\pi,\top} - I)(d^D \cdot \widehat{w}) \rangle| \\ &= \left| \left\langle \widehat{q} - q^{\pi}, (I - \gamma P^{\pi,\top})(d^D \cdot w^{\pi} - d^D \cdot \widehat{w}) \right\rangle \right| \\ &\leq ||(I - \gamma P^{\pi})(\widehat{q} - q^{\pi})||_{2,d^D} ||\widehat{w} - w^{\pi}||_{2,d^D} \end{split}$$

where the last equality is due to the fact that  $(1-\gamma)\mu_0^{\pi} = (I-\gamma P^{\pi})(d^D \cdot w^{\pi})$ , and the final inequality is from Cauchy-Schwarz. We can automatically bound the  $||\hat{w} - w^{\pi}||_{2,d^D}$  term using Theorem 6, and it remains to bound  $||(I - \gamma P^{\pi})(\hat{q} - q^{\pi})||_{2,d^D}$ . We will consider two cases, first when  $d^D > 0$  thus  $\text{Diag}(d^D)$  is invertible, and second, when Assumption 7 is satisfied.

In the first case, let  $D = \text{Diag}(d^D)$ , which by assumption is invertible. Then

$$\begin{aligned} ||(I - \gamma P^{\pi})(\widehat{q} - q^{\pi})||_{2,d^{D}}^{2} &= (\widehat{q} - q^{\pi})^{\top} (I - \gamma P^{\pi})^{\top} D(I - \gamma P^{\pi})(\widehat{q} - q^{\pi}) \\ &= ||D^{1/2}(I - \gamma P^{\pi})(\widehat{q} - q^{\pi})||_{2}^{2} \\ &= ||D^{1/2}(I - \gamma P^{\pi})D^{-1/2}D^{1/2}(\widehat{q} - q^{\pi})||_{2}^{2} \\ &\leq ||D^{1/2}(\widehat{q} - q^{\pi})||_{2}^{2}||D^{1/2}(I - \gamma P^{\pi})D^{-1/2}||_{2}^{2} \\ &= ||\widehat{q} - q^{\pi}||_{d^{D},2}^{2}||I - \gamma P^{\pi}||_{2}^{2} \end{aligned}$$

in the last line using the fact that the eigenvalues of a matrix A and  $L^{-1}AL$  are the same for any invertible matrix L. Thus, denoting the largest singular value of a matrix by  $\sigma_{max}$ ,

$$|\widetilde{J}(\pi) - J(\pi)| \le \sigma_{max} \left( I - \gamma P^{\pi} \right) ||\widehat{w} - w^{\pi}||_{2,d^{D}} ||\widehat{q} - q^{\pi}||_{2,d^{D}}$$

Using Theorem 6 and Theorem 2 in the last line to control the errors of  $\hat{w}$  and  $\hat{q}$  in the last line, followed by a union bound over the three inequalities, gives the result.

For the second case, we can directly apply Lemma 17:

$$\begin{aligned} |J(\pi) - J(\pi)| &\leq ||(I - \gamma P^{\pi})(\widehat{q} - q^{\pi})||_{2,d^{D}} ||\widehat{w} - w^{\pi}||_{2,d^{D}} \\ &\leq \left(||\widehat{q} - q^{\pi}||_{2,d^{D}} + \gamma||P^{\pi}(\widehat{q} - q^{\pi})||_{2,d^{D}}\right) ||\widehat{w} - w^{\pi}||_{2,d^{D}} \\ &\leq \left(1 + \gamma \sqrt{\mathcal{C}_{s'/s}\mathcal{C}_{\pi/\pi^{D}}}\right) ||\widehat{q} - q^{\pi}||_{2,d^{D}} ||\widehat{w} - w^{\pi}||_{2,d^{D}}, \end{aligned}$$

and again applying Theorem 6 and Theorem 2 gives the result.

Lemma 17 uses Assumption 7 to bound the distance in value functions under the transition operator, and is stated and proved below.

Lemma 17. Under Assumption 7,

$$||P^{\pi}(\hat{q}-q^{\pi})||_{2,d^{D}} \leq \sqrt{\mathcal{C}_{s'/s}\mathcal{C}_{\pi/\pi^{D}}}||\hat{q}-q^{\pi}||_{2,d^{D}}.$$

*Proof.* Define  $||P^{\pi}||_{2,d^D} := \sup_{x \neq 0} ||P^{\pi}x||_{2,d^D} / ||x||_{2,d^D}$ . Then

$$||P^{\pi}(\widehat{q}-q^{\pi})||_{2,d^{D}} \leq ||P^{\pi}||_{2,d^{D}} ||\widehat{q}-q^{\pi}||_{2,d^{D}}.$$

It remains to bound  $||P^{\pi}||_{2,d^{D}}$ . For any x,

$$||P^{\pi}x||_{2,d^{D}}^{2} = \mathbb{E}_{(s,a)\sim d^{D}} \left[ \left( \mathbb{E}_{(s',a')\sim P^{\pi}(\cdot|s,a)}[x(s',a')] \right)^{2} \right] \\ \leq \mathbb{E}_{(s,a,s',a')\sim d^{D}\times P^{\pi}}[x(s',a')^{2}] \\ \leq \max_{s,a} \left| \frac{d_{s'}^{D}(s)\pi(a|s)}{d^{D}(s)\pi^{D}(a|s)} \right| \mathbb{E}_{(s,a)\sim d^{D}}[x(s,a)^{2}] \\ = \mathcal{C}_{s'/s}\mathcal{C}_{\pi/\pi^{D}} ||x||_{2,d^{D}}^{2}$$

This implies that  $||P^{\pi}||_{2,d^D} \leq \sqrt{C_{s'/s}C_{\pi/\pi^D}}$ , which gives the stated result.