

---

# CoinDICE: Off-Policy Confidence Interval Estimation

---

\*Bo Dai<sup>1,\*</sup>, Ofir Nachum<sup>1</sup>, Yinlam Chow<sup>1</sup>  
Lihong Li<sup>1</sup>, Csaba Szepesvári<sup>2,3</sup>, Dale Schuurmans<sup>1,3</sup>  
<sup>1</sup>Google Research, Brain Team    <sup>2</sup>University of Alberta    <sup>3</sup>DeepMind

## Abstract

We study *high-confidence behavior-agnostic off-policy evaluation* in reinforcement learning, where the goal is to estimate a confidence interval on a target policy’s value, given only access to a static experience dataset collected by unknown behavior policies. Starting from a function space embedding of the linear program formulation of the  $Q$ -function, we obtain an optimization problem with generalized estimating equation constraints. By applying the generalized empirical likelihood method to the resulting Lagrangian, we propose *CoinDICE*, a novel and efficient algorithm for computing confidence intervals. Theoretically, we prove the obtained confidence intervals are valid, in both asymptotic and finite-sample regimes. Empirically, we show in a variety of benchmarks that the confidence interval estimates are tighter and more accurate than existing methods.<sup>2</sup>

## 1 Introduction

One of the major barriers that hinders the application of reinforcement learning (RL) is the ability to evaluate new policies reliably *before* deployment, a problem generally known as *off-policy evaluation* (OPE). In many real-world domains, *e.g.*, healthcare (Murphy et al., 2001; Gottesman et al., 2018), recommendation (Li et al., 2011; Chen et al., 2019), and education (Mandel et al., 2014), deploying a new policy can be expensive, risky or unsafe. Accordingly, OPE has seen a recent resurgence of research interest, with many methods proposed to estimate the value of a policy (Precup et al., 2000; Dudík et al., 2011; Bottou et al., 2013; Jiang and Li, 2016; Thomas and Brunskill, 2016; Liu et al., 2018; Nachum et al., 2019a; Kallus and Uehara, 2019a,b; Zhang et al., 2020b).

However, the very settings where OPE is necessary usually entail limited data access. In these cases, obtaining knowledge of the uncertainty of the estimate is as important as having a consistent estimator. That is, rather than a *point estimate*, many applications would benefit significantly from having *confidence intervals* on the value of a policy. The problem of estimating these confidence intervals, known as *high-confidence off-policy evaluation* (HCOPE) (Thomas et al., 2015b), is imperative in real-world decision making, where deploying a policy without high-probability safety guarantees can have catastrophic consequences (Thomas, 2015). Most existing high-confidence off-policy evaluation algorithms in RL (Bottou et al., 2013; Thomas et al., 2015a,b; Hanna et al., 2017) construct such intervals using statistical techniques such as concentration inequalities and the bootstrap applied to importance corrected estimates of policy value. The primary challenge with these correction-based approaches is the high variance resulting from multiplying per-step importance ratios in long-horizon problems. Moreover, they typically require full knowledge (or a good estimate) of the behavior policy, which is not easily available in behavior-agnostic OPE settings (Nachum et al., 2019a).

In this work, we propose an algorithm for behavior-agnostic HCOPE. We start from a linear programming formulation of the state-action value function. We show that the value of the policy may be obtained from a Lagrangian optimization problem for generalized estimating equations

---

\*Equal contribution. Email: {bodai, ofirnachum}@google.com.

<sup>2</sup>Open-source code for CoinDICE is available at [https://github.com/google-research/dice\\_rl](https://github.com/google-research/dice_rl).

over data sampled from off-policy distributions. This observation inspires a generalized empirical likelihood approach (Owen, 2001; Broniatowski and Keziou, 2012; Duchi et al., 2016) to confidence interval estimation. These derivations enable us to express high-confidence lower and upper bounds for the policy value as results of minimax optimizations over an arbitrary offline dataset, with the appropriate distribution corrections being implicitly estimated during the optimization. We translate this understanding into a practical estimator, *Confidence Interval Distribution Correction Estimation* (CoinDICE), and design an efficient algorithm for implementing it. We then justify the asymptotic coverage of these bounds and present non-asymptotic guarantees to characterize finite-sample effects. Notably, CoinDICE is behavior-agnostic and its objective function does not involve any per-step importance ratios, and so the estimator is less susceptible to high-variance gradient updates. We evaluate CoinDICE in a number of settings and show that it provides both tighter confidence interval estimates and more correctly matches the desired statistical coverage compared to existing methods.

## 2 Preliminaries

For a set  $W$ , the set of probability measures over  $W$  is denoted by  $\mathcal{P}(W)$ .<sup>3</sup> We consider a Markov Decision Process (MDP) (Puterman, 2014),  $\mathcal{M} = (S, A, T, R, \gamma, \mu_0)$ , where  $S$  denotes the state space,  $A$  denotes the action space,  $T : S \times A \rightarrow \mathcal{P}(S)$  is the transition probability kernel,  $R : S \times A \rightarrow \mathcal{P}([0, R_{\max}])$  is a bounded reward kernel,  $\gamma \in (0, 1]$  is the discount factor, and  $\mu_0$  is the initial state distribution.

A policy,  $\pi : S \rightarrow \mathcal{P}(A)$ , can be used to generate a random trajectory by starting from  $s_0 \sim \mu_0(s)$ , then following  $a_t \sim \pi(s_t)$ ,  $r_t \sim R(s_t, a_t)$  and  $s_{t+1} \sim T(s_t, a_t)$  for  $t \geq 0$ . The state- and action-value functions of  $\pi$  are denoted  $V^\pi$  and  $Q^\pi$ , respectively. The policy also induces an occupancy measure,  $d^\pi(s, a) := (1-\gamma)\mathbb{E}_\pi \left[ \sum_{t \geq 0} \gamma^t \mathbf{1}\{s_t = s, a_t = a\} \right]$ , the normalized discounted probability of visiting  $(s, a)$  in a trajectory generated by  $\pi$ , where  $\mathbf{1}\{\cdot\}$  is the indicator function. Finally, the *policy value* is defined as the *normalized* expected reward accumulated along a trajectory:

$$\rho_\pi := (1-\gamma) \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s_0 \sim \mu_0, a_t \sim \pi(s_t), r_t \sim R(s_t, a_t), s_{t+1} \sim T(s_t, a_t) \right]. \quad (1)$$

We are interested in estimating the policy value and its confidence interval (CI) in the *behavior agnostic off-policy* setting (Nachum et al., 2019a; Zhang et al., 2020a), where interaction with the environment is limited to a static dataset of experience  $\mathcal{D} := \{(s, a, s', r)\}_{i=1}^n$ . Each tuple in  $\mathcal{D}$  is generated according to  $(s, a) \sim d^{\mathcal{D}}, r \sim R(s, a)$  and  $s' \sim T(s, a)$ , where  $d^{\mathcal{D}}$  is an unknown distribution over  $S \times A$ , perhaps induced by one or more unknown behavior policies. The initial distribution  $\mu_0(s)$  is assumed to be easy to sample from, as is typical in practice. Abusing notation, we denote by  $d^{\mathcal{D}}$  both the distribution over  $(s, a, s', r)$  and its marginal on  $(s, a)$ . We use  $\mathbb{E}_d[\cdot]$  for the expectation over a given distribution  $d$ , and  $\mathbb{E}_{\mathcal{D}}[\cdot]$  for its empirical approximation using  $\mathcal{D}$ .

Following previous work (Sutton et al., 2012; Uehara et al., 2019; Zhang et al., 2020a), for ease of exposition we assume the transitions in  $\mathcal{D}$  are *i.i.d.* However, our results may be extended to fast-mixing, ergodic MDPs, where the the empirical distribution of states along a long trajectory is close to being *i.i.d.* (Antos et al., 2008; Lazaric et al., 2012; Dai et al., 2017; Duchi et al., 2016).

Under mild regularity assumptions, the OPE problem may be formulated as a linear program – referred to as the  $Q$ -LP (Nachum et al., 2019b; Nachum and Dai, 2020) – with the following primal and dual forms:

$$\begin{aligned} \min_{Q: S \times A \rightarrow \mathbb{R}} \quad & (1-\gamma) \mathbb{E}_{\mu_0 \pi} [Q(s_0, a_0)] & (2) \\ \text{s.t.} \quad & Q(s, a) \geq R(s, a) + \gamma \cdot \mathcal{P}^\pi Q(s, a), \\ & \forall (s, a) \in S \times A, \end{aligned} \quad \text{and} \quad \begin{aligned} \max_{d: S \times A \rightarrow \mathbb{R}_+} \quad & \mathbb{E}_d [r(s, a)] & (3) \\ \text{s.t.} \quad & d(s, a) = (1-\gamma) \mu_0 \pi(s, a) + \gamma \cdot \mathcal{P}_*^\pi d(s, a), \\ & \forall (s, a) \in S \times A, \end{aligned}$$

where the operator  $\mathcal{P}^\pi$  and its adjoint,  $\mathcal{P}_*^\pi$ , are defined as

$$\begin{aligned} \mathcal{P}^\pi Q(s, a) &:= \mathbb{E}_{s' \sim T(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')], \\ \mathcal{P}_*^\pi d(s, a) &:= \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d(\tilde{s}, \tilde{a}). \end{aligned}$$

<sup>3</sup>All sets and maps are assumed to satisfy appropriate measurability conditions; which we will omit from below for the sake of reducing clutter.

The optimal solutions of (2) and (3) are the  $Q$ -function,  $Q^\pi$ , and stationary state-action occupancy,  $d^\pi$ , respectively, for policy  $\pi$ ; see Nachum et al. (2019b, Theorems 3 & 5) for details as well as extensions to the undiscounted case.

Using the Lagrangian of (2) or (3), we have

$$\rho_\pi = \min_Q \max_{\tau \geq 0} (1 - \gamma) \mathbb{E}_{\mu_0 \pi} [Q(s_0, a_0)] + \mathbb{E}_{d^\pi} [\tau(s, a) (R(s, a) + \gamma Q(s', a') - Q(s, a))], \quad (4)$$

where  $\tau(s, a) := \frac{d(s, a)}{d^\pi(s, a)}$  is the *stationary distribution corrector*. One of the key benefits of the minimax optimization (4) is that both expectations can be immediately approximated by sample averages.<sup>4</sup> In fact, this formulation allows the derivation of several recent behavior-agnostic OPE estimators in a unified manner (Nachum et al., 2019a; Uehara et al., 2019; Zhang et al., 2020a; Nachum and Dai, 2020).

### 3 CoinDICE

We now develop a new approach to obtaining confidence intervals for OPE. The algorithm, *Confidence INterval stationary DIstribution Correction Estimation (CoinDICE)*, is derived by combining function space embedding and the previously described  $Q$ -LP.

#### 3.1 Function Space Embedding of Constraints

Both the primal and dual forms of the  $Q$ -LP contain  $|S| |A|$  constraints that involve expectations over state transition probabilities. Working directly with these constraints quickly becomes computationally and statistically prohibitive when  $|S| |A|$  is large or infinite, as with standard LP approaches (De Farias and Van Roy, 2003). Instead, we consider a relaxation that embeds the constraints in a function space:

$$\tilde{\rho}_\pi := \max_{d: S \times A \rightarrow \mathbb{R}_+} \mathbb{E}_d [r(s, a)] \quad \text{s.t.} \quad \langle \phi, d \rangle = \langle \phi, (1 - \gamma) \mu_0 \pi + \gamma \cdot \mathcal{P}_*^\pi d \rangle, \quad (5)$$

where  $\phi: S \times A \rightarrow \Omega^p \subset \mathbb{R}^p$  is a feature map, and  $\langle \phi, d \rangle := \int \phi(s, a) d(s, a) ds da$ . By projecting the constraints onto a function space with feature mapping  $\phi$ , we can reduce the number of constraints from  $|S| |A|$  to  $p$ . Note that  $p$  may still be infinite. The constraint in (5) can be written as *generalized estimating equations* (Qin and Lawless, 1994; Lam and Zhou, 2017) for the correction ratio  $\tau(s, a)$  over augmented samples  $x := (s_0, a_0, s, a, r, s', a')$  with  $(s_0, a_0) \sim \mu_0 \pi$ ,  $(s, a, r, s') \sim d^\pi$ , and  $a' \sim \pi(\cdot | s')$ ,

$$\langle \phi, d \rangle = \langle \phi, (1 - \gamma) \mu_0 \pi + \gamma \cdot \mathcal{P}_*^\pi d \rangle \Leftrightarrow \mathbb{E}_x [\Delta(x; \tau, \phi)] = 0, \quad (6)$$

where  $\Delta(x; \tau, \phi) := (1 - \gamma) \phi(s_0, a_0) + \tau(s, a) (\gamma \phi(s', a') - \phi(s, a))$ . The corresponding Lagrangian is

$$\tilde{\rho}_\pi = \max_{\tau: S \times A \rightarrow \mathbb{R}_+} \min_{\beta \in \mathbb{R}^p} \mathbb{E}_{d^\pi} [\tau \cdot r(s, a)] + \langle \beta, \mathbb{E}_{d^\pi} [\Delta(x; \tau, \phi)] \rangle. \quad (7)$$

This embedding approach for the dual  $Q$ -LP is closely related to approximation methods for the standard state-value LP (De Farias and Van Roy, 2003; Pazis and Parr, 2011; Lakshminarayanan et al., 2017). The gap between the solutions to (5) and the original dual LP (3) depends on the expressiveness of the feature mapping  $\phi$ . Before stating a theorem that quantifies the error, we first offer a few examples to provide intuition for the role played by  $\phi$ .

**Example (Indicator functions):** Suppose  $p = |S| |A|$  is finite and  $\phi = [\delta_{s,a}]_{(s,a) \in S \times A}$ , where  $\delta_{s,a} \in \{0, 1\}^p$  with  $\delta_{s,a} = 1$  at position  $(s, a)$  and 0 otherwise. Plugging this feature mapping into (5), we recover the original dual  $Q$ -LP (3).

**Example (Full-rank basis):** Suppose  $\Phi \in \mathbb{R}^{p \times p}$  is a full-rank matrix with  $p = |S| |A|$ ; furthermore,  $\phi(s, a) = \Phi((s, a), \cdot)^\top$ . Although the constraints in (5) and (3) are different, their solutions are identical. This can be verified by the Lagrangian in Appendix A.

**Example (RKHS function mappings):** Suppose  $\phi(s, a) := k((s, a), \cdot) \in \mathbb{R}^p$  with  $p = \infty$ , which forms a reproducing kernel Hilbert space (RKHS)  $\mathcal{H}_k$ . The LHS and RHS in the constraint of (5) are the kernel embeddings of  $d(s, a)$  and  $(1 - \gamma) \mu_0 \pi(s, a) + \gamma \cdot \mathcal{P}_*^\pi d(s, a)$  respectively. The constraint in (5) can then be understood as a form of distribution matching by comparing

<sup>4</sup>We assume one can sample initial states from  $\mu_0$ , an assumption that often holds in practice. Then, the data in  $\mathcal{D}$  can be treated as being augmented as  $(s_0, a_0, s, a, r, s', a')$  with  $a_0 \sim \pi(a | s_0)$ ,  $a' \sim \pi(a | s')$ .

kernel embeddings, rather than element-wise matching as in (3). If the kernel function  $k(\cdot, \cdot)$  is characteristic, the embeddings of two distributions will match if and only if the distributions are identical almost surely (Sriperumbudur et al., 2011).

**Theorem 1 (Approximation error)** *Suppose the constant function  $\mathbf{1} \in \mathcal{F}_\phi := \text{span}\{\phi\}$ . Then,*

$$0 \leq \tilde{\rho}_\pi - \rho_\pi \leq 2 \min_{\beta} \|Q^\pi - \langle \beta, \phi \rangle\|_\infty,$$

where  $Q^\pi$  is the fixed-point solution to the Bellman equation  $Q(s, a) = R(s, a) + \gamma \mathcal{P}^\pi Q(s, a)$ .

Please refer to Appendix A for the proof. The condition  $\mathbf{1} \in \mathcal{F}_\phi$  is standard and is trivial to satisfy. Although the approximation error relies on  $\|\cdot\|_\infty$ , a sharper bound that relies on a norm taking the state-action distribution into account can also be obtained (De Farias and Van Roy, 2003). We focus on characterizing the uncertainty due to sampling in this paper, so for ease of exposition we will consider a setting where  $\phi$  is sufficiently expressive to make the approximation error zero. If desired, the approximation error in Theorem 1 can be included in the analysis.

Note that, compared to using a characteristic kernel to ensure injectivity for the RKHS embeddings over all distributions (and thus guaranteeing arbitrarily small approximation error), Theorem 1 only requires that  $Q^\pi$  be represented in  $\mathcal{F}_\phi$ , which is a much weaker condition. In practice, one may also learn the feature mapping  $\phi$  for the projection jointly.

### 3.2 Off-policy Confidence Interval Estimation

By introducing the function space embedding of the constraints in (5), we have transformed the original point-wise constraints in the  $Q$ -LP to generalized estimating equations. This paves the way to applying the generalized empirical likelihood (EL) (Owen, 2001; Broniatowski and Keziou, 2012; Bertail et al., 2014; Duchi et al., 2016) method to estimate a confidence interval on policy value.

Recall that, given a convex, lower-semicontinuous function  $f : \mathbb{R}_+ \rightarrow \mathbb{R}$  satisfying  $f(1) = 0$ , the  $f$ -divergence between densities  $p$  and  $q$  on  $\mathbb{R}$  is defined as  $D_f(P||Q) := \int Q(dx) f\left(\frac{dP(x)}{dQ(x)}\right) dx$ .

Given an  $f$ -divergence, we propose our main confidence interval estimate based on the following confidence set  $C_{n,\xi}^f \subset \mathbb{R}$ :

$$C_{n,\xi}^f := \left\{ \tilde{\rho}_\pi(w) = \max_{\tau \geq 0} \mathbb{E}_w[\tau \cdot r] \mid w \in \mathcal{K}_f, \mathbb{E}_w[\Delta(x; \tau, \phi)] = 0 \right\} \text{ with } \mathcal{K}_f := \left\{ w \in \mathcal{P}^{n-1}(\hat{p}_n), \right. \\ \left. D_f(w||\hat{p}_n) \leq \frac{\xi}{n} \right\}, \quad (8)$$

where  $\mathcal{P}^{n-1}(\hat{p}_n)$  denotes the  $n$ -simplex on the support of  $\hat{p}_n$ , the empirical distribution over  $\mathcal{D}$ . It is easy to verify that this set  $C_{n,\xi}^f \subset \mathbb{R}$  is convex, since  $\tilde{\rho}_\pi(w)$  is a convex function over a convex feasible set. Thus,  $C_{n,\xi}^f$  is an interval. In fact,  $C_{n,\xi}^f$  is the image of the policy value  $\tilde{\rho}_\pi$  on a bounded (in  $f$ -divergence) perturbation to  $w$  in the neighborhood of the empirical distribution  $\hat{p}_n$ .

Intuitively, the confidence interval  $C_{n,\xi}^f$  possesses a close relationship to bootstrap estimators. In vanilla bootstrap, one constructs a set of empirical distributions  $\{w^i\}_{i=1}^m$  by resampling from the dataset  $\mathcal{D}$ . Such subsamples are used to form the empirical distribution on  $\{\tilde{\rho}(w^i)\}_{i=1}^m$ , which provides population statistics for confidence interval estimation. However, this procedure is computationally very expensive, involving  $m$  separate optimizations. By contrast, our proposed estimator  $C_{n,\xi}^f$  exploits the asymptotic properties of the statistic  $\tilde{\rho}_\pi(w)$  to derive a target confidence interval by solving only *two* optimization problems (Section 3.3), a dramatic savings in computational cost.

Before introducing the algorithm for computing  $C_{n,\xi}^f$ , we establish the first key result that, by choosing  $\xi = \chi_{(1)}^{2,1-\alpha}$ ,  $C_{n,\xi}^f$  is asymptotically a  $(1 - \alpha)$ -confidence interval on the policy value, where  $\chi_{(1)}^{2,1-\alpha}$  is the  $(1 - \alpha)$ -quantile of the  $\chi^2$ -distribution with 1 degree of freedom.

**Theorem 2 (Informal asymptotic coverage)** *Under some mild conditions, if  $\mathcal{D}$  contains i.i.d. samples and the optimal solution to the Lagrangian of (5) is unique, we have*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\rho_\pi \in C_{n,\xi}^f\right) = \mathbb{P}\left(\chi_{(1)}^2 \leq \xi\right). \quad (9)$$

Thus,  $C_{n,\chi_{(1)}^2,1-\alpha}^f$  is an asymptotic  $(1 - \alpha)$ -confidence interval of the value of the policy  $\pi$ .

Please refer to Appendix E.1 for the precise statement and proof of Theorem 2.

Theorem 2 generalizes the result in Duchi et al. (2016) to statistics with generalized estimating equations, maintaining the 1 degree of freedom in the asymptotic  $\chi_{(1)}^2$ -distribution. One may also apply existing results for EL with generalized estimating equations (e.g., Lam and Zhou, 2017), but these would lead to a limiting distribution of  $\chi_{(m)}^2$  with  $m \gg 1$  degrees of freedom, resulting in a much looser confidence interval estimate than Theorem 2.

Note that Theorem 2 can also be specialized to multi-armed contextual bandits to achieve a tighter confidence interval estimate in this special case. In particular, for contextual bandits, the stationary distribution constraint in (5),  $\mathbb{E}_w [\Delta(x; \tau, \phi)] = 0$ , is no longer needed, and can be replaced by  $\mathbb{E}_w [\tau - 1] = 0$ . Then by the same technique used for MDPs, we can obtain a confidence interval estimate for offline contextual bandits; see details in Appendix C. Interestingly, the resulting confidence interval estimate not only has the same asymptotic coverage as previous work (Karampatziakis et al., 2019), but is also simpler and computationally more efficient.

### 3.3 Computing the Confidence Interval

Now we provide a distributional robust optimization view of the upper and lower bounds of  $C_{n,\xi}^f$ .

**Theorem 3 (Upper and lower confidence bounds)** Denote the upper and lower confidence bounds of  $C_{n,\xi}^f$  by  $u_n$  and  $l_n$ , respectively:

$$[l_n, u_n] = \left[ \min_{w \in \mathcal{K}_f} \min_{\beta \in \mathbb{R}^p} \max_{\tau \geq 0} \mathbb{E}_w [\ell(x; \tau, \beta)], \max_{w \in \mathcal{K}_f} \max_{\tau \geq 0} \min_{\beta \in \mathbb{R}^p} \mathbb{E}_w [\ell(x; \tau, \beta)] \right], \quad (10)$$

$$= \left[ \min_{\beta \in \mathbb{R}^p} \max_{\tau \geq 0} \min_{w \in \mathcal{K}_f} \mathbb{E}_w [\ell(x; \tau, \beta)], \max_{\tau \geq 0} \min_{\beta \in \mathbb{R}^p} \max_{w \in \mathcal{K}_f} \mathbb{E}_w [\ell(x; \tau, \beta)] \right], \quad (11)$$

where  $\ell(x; \tau, \beta) := \tau \cdot r + \beta^\top \Delta(x; \tau, \phi)$ . For any  $(\tau, \beta, \lambda, \eta)$  that satisfies the constraints in (11), the optimal weights for the upper and lower confidence bounds are

$$w_l = f_*' \left( \frac{\eta - \ell(x; \tau, \beta)}{\lambda} \right) \quad \text{and} \quad w_u = f_*' \left( \frac{\ell(x; \tau, \beta) - \eta}{\lambda} \right). \quad (12)$$

respectively. Therefore, the confidence bounds can be simplified as:

$$\begin{bmatrix} l_n \\ u_n \end{bmatrix} = \begin{bmatrix} \min_{\beta} \max_{\tau \geq 0, \lambda \geq 0, \eta} \mathbb{E}_{\mathcal{D}} \left[ -\lambda f_*' \left( \frac{\eta - \ell(x; \tau, \beta)}{\lambda} \right) + \eta - \lambda \frac{\xi}{n} \right] \\ \max_{\tau \geq 0} \min_{\beta, \lambda \geq 0, \eta} \mathbb{E}_{\mathcal{D}} \left[ \lambda f_*' \left( \frac{\ell(x; \tau, \beta) - \eta}{\lambda} \right) + \eta + \lambda \frac{\xi}{n} \right] \end{bmatrix}. \quad (13)$$

The proof of this result relies on Lagrangian duality and the convexity and concavity of the optimization; it may be found in full detail in Appendix D.1.

As we can see in Theorem 3, by exploiting strong duality properties to move  $w$  into the inner most optimizations in (11), the obtained optimization (11) is the distributional robust optimization extension of the saddle-point problem. The closed-form reweighting scheme is demonstrated in (12). For particular  $f$ -divergences, such as the  $KL$ - and 2-power divergences, for a fixed  $(\beta, \tau)$ , the optimal  $\eta$  can be easily computed and the weights  $w$  recovered in closed-form. For example, by using  $KL(w || \hat{p}_n)$ , (12) can be used to obtain the updates

$$w_l(x) = \exp \left( \frac{\eta_l - \ell(x; \tau, \beta)}{\lambda} \right), \quad w_u(x) = \exp \left( \frac{\ell(x; \tau, \beta) - \eta_u}{\lambda} \right), \quad (14)$$

where  $\eta_l$  and  $\eta_u$  provide the normalizing constants. (For closed-form updates of  $w$  w.r.t. other  $f$ -divergences, please refer to Appendix D.2.) Plug the closed-form of optimal weights into (11), this greatly simplifies the optimization over the data perturbations yielding (13), and establishes the connection to the prioritized experiences replay (Schaul et al., 2016), where both reweight the experience data according to their loss, but with different reweighting schemes.

Note that it is straightforward to check that the estimator for  $u_n$  in (13) is nonconvex-concave and the estimator for  $l_n$  in (13) is nonconcave-convex. Therefore, one could alternatively apply stochastic gradient descent-ascent (SGDA) for to solve (13) and benefit from attractive finite-step convergence guarantees (Lin et al., 2019).



**Remark (Practical considerations):** As also observed in [Namkoong and Duchi \(2016\)](#), SGDA for (13) could potentially suffer from high variance in both the objective and gradients when  $\lambda$  approaches 0. In [Appendix D.3](#), we exploit several properties of (11), which leads to a computational efficient algorithm, to overcome the numerical issue. Please refer to [Appendix D.3](#) for the details of [Algorithm 1](#) and the practical considerations.

**Remark (Joint learning for feature embeddings):** The proposed framework also allows for the possibility to learn the features for constraint projection. In particular, consider  $\zeta(\cdot, \cdot) := \beta^\top \phi(\cdot, \cdot) : S \times A \rightarrow \mathbb{R}$ . Note that we could treat the combination  $\beta^\top \phi(s, a)$  together as the Lagrange multiplier function for the original  $Q$ -LP with *infinitely* many constraints, hence both  $\beta$  and  $\phi(\cdot, \cdot)$  could be updated jointly. Although the conditions for asymptotic coverage no longer hold, the finite-sample correction results of the next section are still applicable. This might offer an interesting way to reduce the approximation error introduced by inappropriate feature embeddings of the constraints, while still maintaining calibrated confidence intervals.

## 4 Finite-sample Analysis

[Theorem 2](#) establishes the asymptotic  $(1 - \alpha)$ -coverage of the confidence interval estimates produced by CoinDICE, ignoring higher-order error terms that vanish as sample size  $n \rightarrow \infty$ . In practice, however,  $n$  is always finite, so it is important to quantify these higher-order terms. This section addresses this problem, and presents a finite-sample bound for the estimate of CoinDICE. In the following, we let  $\mathcal{F}_\tau$  and  $\mathcal{F}_\beta$  be the function classes of  $\tau$  and  $\beta$  used by CoinDICE.

**Theorem 4 (Informal finite-sample correction)** *Denote by  $d_{\mathcal{F}_\tau}$  and  $d_{\mathcal{F}_\beta}$  the finite VC-dimension of  $\mathcal{F}_\tau$  and  $\mathcal{F}_\beta$ , respectively. Under some mild conditions, when  $D_f$  is  $\chi^2$ -divergence, we have*

$$\mathbb{P}(\rho_\pi \in [l_n - \kappa_n, u_n + \kappa_n]) \geq 1 - 12 \exp\left(c_1 + 2(d_{\mathcal{F}_\tau} + d_{\mathcal{F}_\beta} - 1) \log n - \frac{\xi}{18}\right),$$

where  $c_1 = 2c + \log d_{\mathcal{F}_\tau} + \log d_{\mathcal{F}_\beta} + (d_{\mathcal{F}_\tau} + d_{\mathcal{F}_\beta} - 1)$ ,  $\kappa_n = \frac{11M\xi}{6n} + 2\frac{C_\ell M}{n} \left(1 + 2\sqrt{\frac{\xi}{9n}}\right)$ , and  $(c, M, C_\ell)$  are universal constants.

The precise statement and detailed proof of [Theorem 4](#) can be found in [Appendix E.2](#). The proof relies on empirical Bernstein bounds with a careful analysis of the variance term. Compared to the vanilla sample complexity of  $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ , we achieve a faster rate of  $\mathcal{O}\left(\frac{1}{n}\right)$  without any additional assumptions on the noise or curvature conditions. The tight sample complexity in [Theorem 4](#) implies that one can construct the  $(1 - \alpha)$ -finite sample confidence interval by optimizing (11) with  $\xi = 18 \left(\log \frac{\alpha}{12} - c_1 - 2(d_{\mathcal{F}_\tau} + d_{\mathcal{F}_\beta} - 1) \log n\right)$ , and composing with  $\kappa_n$ . However, we observe that this bound can be conservative compared to the asymptotic confidence interval in [Theorem 2](#). Therefore, we will evaluate the asymptotic version of CoinDICE based on [Theorem 2](#) in the experiment.

The conservativeness arises from the use of a union bound. However, we conjecture that the rate is optimal up to a constant. We exploit the VC dimension due to its generality. In fact, the bound can be improved by considering a data-dependent measure, *e.g.*, Rademacher complexity, or by some function class dependent measure, *e.g.*, function norm in RKHS, for specific function approximators.

## 5 Optimism vs. Pessimism Principle

CoinDICE provide both upper and lower bounds of the target policy’s estimated value, which paves the path for applying the principle of optimism ([Lattimore and Szepesvári, 2020](#)) or pessimism ([Swaminathan and Joachims, 2015](#)) in the face of uncertainty for policy optimization in different learning settings.

**Optimism in the face of uncertainty.** Optimism in the face of uncertainty leads to *risk-seeking* algorithms, which can be used to balance the exploration/exploitation trade-off. Conceptually, they always treat the environment as the best plausibly possible. This principle has been successfully applied to stochastic bandit problems, leading to many instantiations of UCB algorithms ([Lattimore and Szepesvári, 2020](#)). In each round, an action is selected according to the upper confidence bound, and the obtained reward will be used to refine the confidence bound iteratively. When applied to MDPs, this principle inspires many optimistic model-based ([Bartlett and Mendelson, 2002](#); [Auer](#)

et al., 2009; Strehl et al., 2009; Szita and Szepesvari, 2010; Dann et al., 2017), value-based (Jin et al., 2018), and policy-based algorithms (Cai et al., 2019). Most of these algorithms are not compatible with function approximators.

We can also implement the optimism principle by optimizing the upper bound in CoinDICE iteratively, *i.e.*,  $\max_{\pi} u_{\mathcal{D}}(\pi)$ . In  $t$ -th iteration, we calculate the gradient of  $u_{\mathcal{D}}(\pi^t)$ , *i.e.*,  $\nabla_{\pi} u_{\mathcal{D}}(\pi^t)$ , based on the existing dataset  $\mathcal{D}_t$ , then, the policy  $\pi_t$  will be updated by (natural) policy gradient and samples will be collected through the updated policy  $\pi_{t+1}$ . Please refer to Appendix F for the gradient computation and algorithm details.

**Pessimism in the face of uncertainty.** In offline reinforcement learning (Lange et al., 2012; Fujimoto et al., 2019; Wu et al., 2019; Nachum et al., 2019b), only a fixed set of data from behavior policies is given, a safe optimization criterion is to maximize the worst-case performance among a set of statistically plausible models (Laroche et al., 2019; Kumar et al., 2019; Yu et al., 2020). In contrast to the previous case of online exploration, this is a pessimism principle (Cohen and Hutter, 2020; Buckman et al., 2020) or counterfactual risk minimization (Swaminathan and Joachims, 2015), and highly related to robust MDP (Iyengar, 2005; Nilim and El Ghaoui, 2005; Tamar et al., 2013; Chow et al., 2015).

Different from most of the existing methods where the worst-case performance is characterized by model-based perturbation or ensemble, the proposed CoinDICE provides a lower bound to implement the pessimism principle, *i.e.*,  $\max_{\pi} l_{\mathcal{D}}(\pi)$ . Conceptually, we apply the (natural) policy gradient w.r.t.  $l_{\mathcal{D}}(\pi^t)$  to update the policy iteratively. Since we are dealing with policy optimization in the offline setting, the dataset  $\mathcal{D}$  keeps unchanged. Please refer to Appendix F for the algorithm details.

## 6 Related Work

Off-policy estimation has been extensively studied in the literature, given its practical importance. Most existing methods are based on the core idea of importance reweighting to correct for distribution mismatches between the target policy and the off-policy data (Precup et al., 2000; Bottou et al., 2013; Li et al., 2015; Xie et al., 2019). Unfortunately, when applied naively, importance reweighting can result in an excessively high variance, which is known as the “curse of horizon” (Liu et al., 2018). To avoid this drawback, there has been rapidly growing interest in estimating the correction ratio of the *stationary* distribution (e.g., Liu et al., 2018; Nachum et al., 2019a; Uehara et al., 2019; Liu et al., 2019; Zhang et al., 2020a,b). This work is along the same line and thus applicable in long-horizon problems. Other off-policy approaches are also possible, notably model-based (e.g., Fonteneau et al., 2013) and doubly robust methods (Jiang and Li, 2016; Thomas and Brunskill, 2016; Tang et al., 2020; Uehara et al., 2019). These techniques can potentially be combined with our algorithm, which we leave for future investigation.

While most OPE works focus on obtaining accurate *point* estimates, several authors provide ways to quantify the amount of uncertainty in the OPE estimates. In particular, confidence bounds have been developed using the central limit theorem (Bottou et al., 2013), concentration inequalities (Thomas et al., 2015b; Kuzborskij et al., 2020), and nonparametric methods such as the bootstrap (Thomas et al., 2015a; Hanna et al., 2017). In contrast to these works, the CoinDICE is asymptotically pivotal, meaning that there are no hidden quantities we need to estimate, which is based on correcting for the stationary distribution in the *behavior-agnostic* setting, thus avoiding the curse of horizon and broadening the application of the uncertainty estimator. Recently, Jiang and Huang (2020) provide confidence intervals for OPE, but focus on the intervals determined by the *approximation error* induced by a function approximator, while our confidence intervals quantify *statistical error*.

Empirical likelihood (Owen, 2001) is a powerful tool with many applications in statistical inference like econometrics (Chen et al., 2018), and more recently in distributionally robust optimization (Duchi et al., 2016; Lam and Zhou, 2017). EL-based confidence intervals can be used to guide exploration in multiarmed bandits (Honda and Takemura, 2010; Cappé et al., 2013), and for OPE (Karampatziakis et al., 2019; Kallus and Uehara, 2019b). While the work of Kallus and Uehara (2019b) is also based on EL, it differs from the present work in two important ways. First, their focus is on developing an asymptotically efficient OPE *point* estimate, not confidence intervals. Second, they solve for timestep-dependent weights, whereas we only need to solve for timestep-*independent* weights from a system of moment matching equations induced by an underlying ergodic Markov chain.

## 7 Experiments

We now evaluate the empirical performance of CoinDICE, comparing it to a number of existing confidence interval estimators for OPE based on concentration inequalities. Specifically, given a dataset of logged trajectories, we first use weighted step-wise importance sampling (Precup et al., 2000) to calculate a separate estimate of the target policy value for each trajectory. Then given such a finite sample of estimates, we then use the empirical *Bernstein* inequality (Thomas et al., 2015b) to derive high-confidence lower and upper bounds for the true value. Alternatively, one may also use *Student's t-test* or Efron's bias corrected and accelerated *bootstrap* (Thomas et al., 2015a).

We begin with a simple bandit setting, devising a two-armed bandit problem with stochastic payoffs. We define the target policy as a near-optimal policy, which chooses the optimal arm with probability 0.95. We collect off-policy data using a behavior policy which chooses the optimal arm with probability of only 0.55. Our results are presented in Figure 1. We plot the empirical coverage and width of the estimated intervals across different confidence levels. More specifically, each data point in Figure 1 is the result of 200 experiments. In each experiment, we randomly sample a dataset and then compute a confidence interval. The *interval coverage* is then computed as the proportion of intervals out of 200 that contain the true value of the target policy. The *interval log-width* is the median of the log of the width of the 200 computed intervals. Figure 1 shows that the intervals produced by CoinDICE achieve an empirical coverage close to the intended coverage. In this simple bandit setting, the coverages of Student's *t* and bootstrapping are also close to correct, although they suffer more in the low-data regime. Notably, the width of the intervals produced by CoinDICE are especially narrow while maintaining accurate coverage.

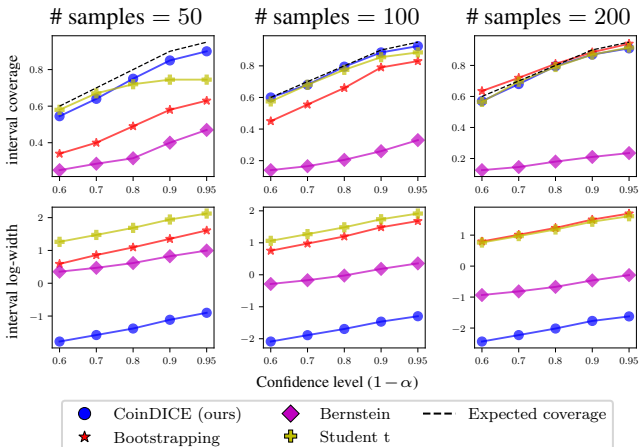


Figure 1: Results of CoinDICE and baseline methods on a simple two-armed bandit. We plot empirical coverage and median log-width (*y*-axes) of intervals evaluated at a number of desired confidence levels (*x*-axis), as measured over 200 random trials. We find that CoinDICE achieves more accurate coverage and narrower intervals compared to the baseline confidence interval estimation methods.

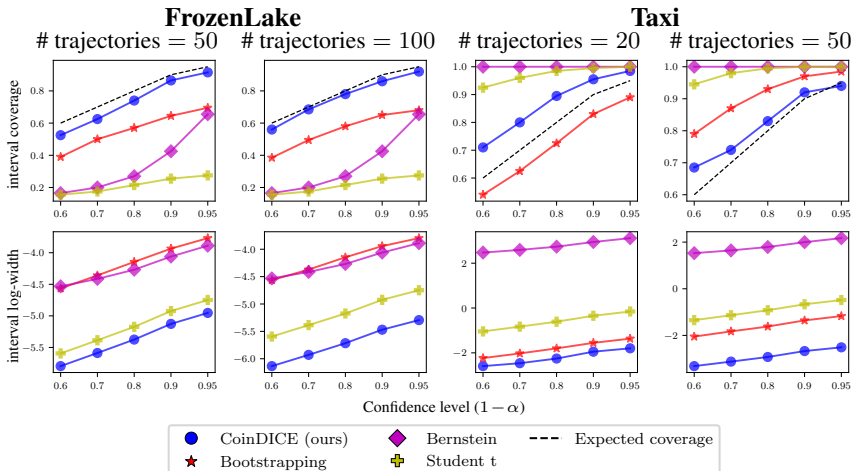


Figure 2: Results of CoinDICE and baseline methods on an infinite-horizon version of FrozenLake and Taxi. In FrozenLake, each dataset consists of trajectories of length 100; in Taxi, each dataset consists of trajectories of length 500.



We now turn to more complicated MDP environments. We use FrozenLake (Brockman et al., 2016), a highly stochastic gridworld environment, and Taxi (Dietterich, 1998), an environment with a moderate state space of 2000 elements. As in (Liu et al., 2018), we modify these environments to be infinite horizon by randomly resetting the state upon termination. The discount factor is  $\gamma = 0.99$ . The target policy is taken to be a near-optimal one, while the behavior policy is highly suboptimal. The behavior policy in FrozenLake is the optimal policy with 0.2 white noise, which reduces the policy value dramatically, from 0.74 to 0.24. For the behavior policies in Taxi and Reacher, we follow the same experiment setting for constructing the behavior policies to collect data as in (Nachum et al., 2019a; Liu et al., 2018).

We follow the same evaluation protocol as in the bandit setting, measuring empirical interval coverage and log-width over 200 experimental trials for various dataset sizes and confidence levels. Results are shown in Figure 2. We find a similar conclusion that CoinDICE consistently achieves more accurate coverage and smaller widths than baselines. Notably, the baseline methods’ accuracy suffers more significantly compared to the simpler bandit setting described earlier.

Lastly, we evaluate CoinDICE on Reacher (Brockman et al., 2016; Todorov et al., 2012), a continuous control environment. In this setting, we use a one-hidden-layer neural network with ReLU activations. Results are shown in Figure 3. To account for the approximation error of the used neural network, we measure the coverage of CoinDICE with respect to a true value computed as the median of a large ensemble of neural networks trained on the off-policy data. To keep the comparison fair, we measure the coverage of the IS-based baselines with respect to a true value computed as the median of a large number of IS-based point estimates. The results show similar conclusions as before: CoinDICE achieves more accurate coverage than the IS-based methods. Still, we see that CoinDICE coverage suffers in this regime, likely due to optimization difficulties. If the optimum of the Lagrangian is only approximately found, the empirical coverage will inevitably be inexact.

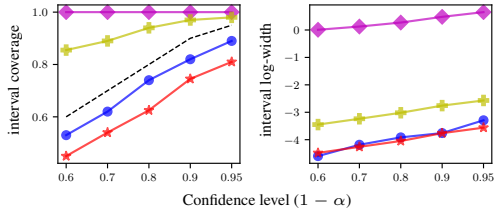


Figure 3: Results of CoinDICE and baseline methods on Reacher (Brockman et al., 2016; Todorov et al., 2012), using 25 trajectories of length 100. Colors and markers are as defined in the legends of previous figures.

## 8 Conclusion

In this paper, we have developed CoinDICE, a novel and efficient confidence interval estimator applicable to the *behavior-agnostic offline* setting. The algorithm builds on a few technical components, including a new feature embedded  $Q$ -LP, and a generalized empirical likelihood approach to confidence interval estimation. We analyzed the asymptotic coverage of CoinDICE’s estimate, and provided an inite-sample bound. On a variety of off-policy benchmarks we empirically compared the new algorithm with several strong baselines and found it to be superior to them.

## Broader Impact

This research is fundamental and targets a broad question in reinforcement learning. The ability to reliably assess uncertainty in off-policy evaluation would have significant positive benefits for safety-critical applications of reinforcement learning. Inaccurate uncertainty estimates create the danger of misleading decision makers and could lead to detrimental consequences. However, our primary goal is to improve these estimators and reduce the ultimate risk of deploying reinforcement-learned systems. The techniques are general and do not otherwise target any specific application area.

## Acknowledgements

We thank Hanjun Dai, Mengjiao Yang and other members of the Google Brain team for helpful discussions. Csaba Szepesvári gratefully acknowledges funding from the Canada CIFAR AI Chairs Program, Amii and NSERC.

## References

- András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pages 8139–8148, 2019.
- P. Auer. Using upper confidence bounds for online learning. In *Proc. 41st Annual Symposium on Foundations of Computer Science*, pages 270–279. IEEE Computer Society Press, Los Alamitos, CA, 2000.
- Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. In *Advances in neural information processing systems*, pages 89–96, 2009.
- P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Patrice Bertail, Emmanuelle Gautherat, and Hugo Harari-Kermadec. Empirical  $\varphi_*$ -divergence minimizers for Hadamard differentiable functionals. In *Topics in Nonparametric Statistics*, pages 21–32. Springer, 2014.
- Léon Bottou, Jonas Peters, Joaquin Quiñero-Candela, Denis X Charles, D Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning systems: The example of computational advertising. *The Journal of Machine Learning Research*, 14(1):3207–3260, 2013.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Michel Broniatowski and Amor Keziou. Divergences and duality for estimation and test under moment condition models. *Journal of Statistical Planning and Inference*, 142(9):2554–2573, 2012.
- Jacob Buckman, Carles Gelada, and Marc G Bellemare. The importance of pessimism in fixed-dataset policy optimization. *arXiv preprint arXiv:2009.06799*, 2020.
- Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*, 2019.
- Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Maillard, Rémi Munos, Gilles Stoltz, et al. Kullback-Leibler upper confidence bounds for optimal sequential allocation. *The Annals of Statistics*, 41(3):1516–1541, 2013.
- Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. Top- $k$  off-policy correction for a REINFORCE recommender system. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pages 456–464, 2019.
- Xiaohong Chen, Timothy M. Christensen, and Elie Tamer. Monte Carlo confidence sets for identified sets. *Econometrica*, 86(6):1965–2018, 2018.
- Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-making: a cvar optimization approach. In *Advances in Neural Information Processing Systems*, pages 1522–1530, 2015.
- Michael K Cohen and Marcus Hutter. Pessimism about unknown unknowns inspires conservatism. In *Conference on Learning Theory*, pages 1344–1373. PMLR, 2020.
- Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. SBEED: Convergent reinforcement learning with nonlinear function approximation. *CoRR*, abs/1712.10285, 2017.

- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 5713–5723, 2017.
- Daniela Pucci De Farias and Benjamin Van Roy. The linear programming approach to approximate dynamic programming. *Operations research*, 51(6):850–865, 2003.
- Thomas G. Dietterich. The MAXQ method for hierarchical reinforcement learning. In *Proc. Intl. Conf. Machine Learning*, pages 118–126. Morgan Kaufmann, San Francisco, CA, 1998.
- John Duchi, Peter Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. *arXiv preprint arXiv:1610.03425*, 2016.
- Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. In *Proceedings of the 28th International Conference on Machine Learning*, pages 1097–1104, 2011. CoRR abs/1103.4601.
- Ivar Ekeland and Roger Temam. *Convex analysis and variational problems*, volume 28. SIAM, 1999.
- Raphael Fonteneau, Susan A. Murphy, Louis Wehenkel, and Damien Ernst. Batch mode reinforcement learning based on the synthesis of artificial trajectories. *Annals of Operations Research*, 208(1): 383–416, 2013.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pages 2052–2062, 2019.
- Omer Gottesman, Fredrik Johansson, Joshua Meier, Jack Dent, Donghun Lee, Srivatsan Srinivasan, Linying Zhang, Yi Ding, David Wihl, Xuefeng Peng, Jiayu Yao, Isaac Lage, Christopher Mosch, Li wei H. Lehman, Matthieu Komorowski, Matthieu Komorowski, Aldo Faisal, Leo Anthony Celi, David Sontag, and Finale Doshi-Velez. Evaluating reinforcement learning algorithms in observational health settings, 2018. arXiv:1805.12298.
- Josiah P. Hanna, Peter Stone, and Scott Niekum. Bootstrapping with models: Confidence intervals for off-policy evaluation. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, pages 4933–4934, 2017.
- Junya Honda and Akimichi Takemura. An asymptotically optimal bandit algorithm for bounded support models. In *COLT*, pages 67–79, 2010.
- Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2): 257–280, 2005.
- Nan Jiang and Jiawei Huang. Minimax confidence interval for off-policy evaluation and policy optimization, 2020. arXiv:2002.02081.
- Nan Jiang and Lihong Li. Doubly robust off-policy evaluation for reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 652–661, 2016.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient? In *Advances in Neural Information Processing Systems*, pages 4863–4873, 2018.
- Nathan Kallus and Masatoshi Uehara. Double reinforcement learning for efficient off-policy evaluation in Markov decision processes. *arXiv preprint arXiv:1908.08526*, 2019a.
- Nathan Kallus and Masatoshi Uehara. Intrinsically efficient, stable, and bounded off-policy evaluation for reinforcement learning. In *Advances in Neural Information Processing Systems 32*, pages 3320–3329, 2019b.
- Nikos Karampatziakis, John Langford, and Paul Mineiro. Empirical likelihood for contextual bandits. *arXiv preprint arXiv:1906.03323*, 2019.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy Q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing Systems*, pages 11784–11794, 2019.

- Ilja Kuzborskij, Claire Vernade, András György, Csaba Szepesvári. Confident off-policy evaluation and selection through self-normalized importance weighting. *arXiv preprint arXiv:2006.10460*, 2020.
- Chandrashekar Lakshminarayanan, Shalabh Bhatnagar, and Csaba Szepesvari. A linearly relaxed approximate linear program for Markov decision processes. *arXiv preprint arXiv:1704.02544*, 2017.
- Henry Lam and Enlu Zhou. The empirical likelihood approach to quantifying uncertainty in sample average approximation. *Operations Research Letters*, 45(4):301–307, 2017.
- Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In *Reinforcement learning*, pages 45–73. Springer, 2012.
- Romain Laroche, Paul Trichelair, and Remi Tachet Des Combes. Safe policy improvement with baseline bootstrapping. In *International Conference on Machine Learning*, pages 3652–3661. PMLR, 2019.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of least-squares policy iteration. *Journal of Machine Learning Research*, 13(Oct):3041–3074, 2012.
- Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-bandit-based news article recommendation algorithms. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 297–306. ACM, 2011.
- Lihong Li, Rémi Munos, and Csaba Szepesvári. Toward minimax off-policy value estimation. pages 608–616, 2015.
- Tianyi Lin, Chi Jin, and Michael I. Jordan. On gradient descent ascent for nonconvex-concave minimax problems. *CoRR*, abs/1906.00331, 2019.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 5356–5366. Curran Associates, Inc., 2018.
- Yao Liu, Pierre-Luc Bacon, and Emma Brunskill. Understanding the curse of horizon in off-policy evaluation via conditional importance sampling, 2019. *arXiv:1910.06508*.
- Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. 2014.
- Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*, 2009.
- Susan A Murphy, Mark J van der Laan, James M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423, 2001.
- Ofir Nachum and Bo Dai. Reinforcement learning via Fenchel-Rockafellar duality. *arXiv preprint arXiv:2001.01866*, 2020.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections. pages 2315–2325, 2019a.
- Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. AlgaeDICE: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019b.
- Hongseok Namkoong and John C Duchi. Stochastic gradient methods for distributionally robust optimization with f-divergences. In *Advances in neural information processing systems*, pages 2208–2216, 2016.

- Hongseok Namkoong and John C Duchi. Variance-based regularization with convex objectives. In *Advances in neural information processing systems*, pages 2971–2980, 2017.
- Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- Art B Owen. *Empirical likelihood*. Chapman and Hall/CRC, 2001.
- Jason Papis and Ronald Parr. Non-parametric approximate linear programming for MDPs. In *AAAI*, 2011.
- Doina Precup, R. S. Sutton, and S. Singh. Eligibility traces for off-policy policy evaluation. In *Proc. Intl. Conf. Machine Learning*, pages 759–766. Morgan Kaufmann, San Francisco, CA, 2000.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Jin Qin and Jerry Lawless. Empirical likelihood and general estimating equations. *the Annals of Statistics*, pages 300–325, 1994.
- R Tyrrell Rockafellar. Augmented lagrange multiplier functions and duality in nonconvex programming. *SIAM Journal on Control*, 12(2):268–285, 1974.
- Werner Römisch. Delta method, infinite dimensional. *Wiley StatsRef: Statistics Reference Online*, 2014.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In *Proceedings of the 4th International Conference on Learning Representations*, 2016.
- B. Sriperumbudur, K. Fukumizu, and G. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12:2389–2410, 2011.
- Alexander L. Strehl, Lihong Li, and Michael L. Littman. Reinforcement learning in finite MDPs: PAC analysis. In *Journal of Machine Learning Research*, 10:2413–2444, 2009.
- Richard S Sutton, Csaba Szepesvári, Alborz Geramifard, and Michael P Bowling. Dyna-style planning with linear function approximation and prioritized sweeping. *arXiv preprint arXiv:1206.3285*, 2012.
- Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *International Conference on Machine Learning*, pages 814–823, 2015.
- Istvan Szita and Csaba Szepesvari. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, page 1031–1038. Omnipress, 2010.
- Aviv Tamar, Huan Xu, and Shie Mannor. Scaling up robust MDPs by reinforcement learning. *arXiv preprint arXiv:1306.6189*, 2013.
- Ziyang Tang, Yihao Feng, Lihong Li, Dengyong Zhou, and Qiang Liu. Doubly robust bias reduction in infinite horizon off-policy estimation. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.
- Philip Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High confidence policy improvement. In *International Conference on Machine Learning*, pages 2380–2388, 2015a.
- Philip S Thomas. *Safe reinforcement learning*. PhD thesis, University of Massachusetts Libraries, 2015.
- Philip S. Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 2139–2148, 2016.
- Philip S Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High-confidence off-policy evaluation. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015b.



- Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 5026–5033. IEEE, 2012.
- Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and Q-function learning for off-policy evaluation. *arXiv preprint arXiv:1910.12809*, 2019.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes*. Springer, 1996.
- Weiran Wang and Miguel A Carreira-Perpinán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. *arXiv preprint arXiv:1309.1541*, 2013.
- Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. In *Advances in Neural Information Processing Systems 32*, pages 9665–9675, 2019.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. MOPE: Model-based offline policy optimization. *arXiv preprint arXiv:2005.13239*, 2020.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. GenDICE: Generalized offline estimation of stationary values. In *International Conference on Learning Representations*, 2020a.
- Shangdong Zhang, Bo Liu, and Shimon Whiteson. GradientDICE: Rethinking generalized offline estimation of stationary values, 2020b. *arXiv:2001.11113*.