
Minimax-Optimal Off-Policy Evaluation with Linear Function Approximation

Yaqi Duan¹ Zeyu Jia² Mengdi Wang^{3,4}

Abstract

This paper studies the statistical theory of off-policy policy evaluation with function approximation in batch data reinforcement learning problem. We consider a regression-based fitted Q iteration method, and show that it is equivalent to a model-based method that estimates a conditional mean embedding of the transition operator. We prove that this method is information-theoretically optimal and has nearly minimal estimation error. In particular, by leveraging contraction property of Markov processes and martingale concentration, we establish a finite-sample instance-dependent error upper bound and a nearly-matching minimax lower bound. The policy evaluation error depends sharply on a restricted χ^2 -divergence over the function class between the long-term distribution of target policy and the distribution of past data. This restricted χ^2 -divergence characterizes the statistical limit of off-policy evaluation, and is both instance-dependent and function-class-dependent. Further, we provide an easily computable confidence bound for the policy evaluator, which may be useful for optimistic planning and safe policy improvement.

1. Introduction

Batch data reinforcement learning (RL) is common in decision-making applications where rich experiences are available but new experiments are costly. A first-order question is how much one can learn from existing experiences to predict and improve the performance of new policies. This is known as the off-policy policy evaluation (OPE) problem, where one needs to estimate the cumulative rewards (aka

^{*}Equal contribution ¹Department of Operations Research and Financial Engineering, Princeton University, NJ, USA ²School of Mathematics, Peking University, Beijing, China ³Department of Operations Research and Financial Engineering, Princeton University, NJ, USA ⁴DeepMind, London, UK. Correspondence to: Mengdi Wang <mengdiw@princeton.edu>.

value) to be earned by a new policy based on logged history.

In this paper, we study the off-policy evaluation using linear function approximation. We assume that the Q-functions of interests belong to a known function class \mathcal{Q} with d basis functions. We adopt a direct regression-based approach and investigate the basic fitted Q iteration (FQI) (Bertsekas et al., 1995; Sutton & Barto, 2018). It works by iteratively estimating Q-functions via supervised learning using the batch data. This approach turns out to be equivalent to the model-based plug-in estimator where one estimates the conditional mean embedding of the unknown transition model and uses it to compute a plug-in value estimator. It is also related to variants of importance sampling methods (see discussions in Sections 1.1 and 3.3).

We provide a finite-sample error upper bound for this policy evaluator, as well as a nearly matching minimax-optimal lower bound. Putting them together, we see that the regression-based policy evaluator is nearly statistically optimal. For RL with horizon H , the minimax-optimal OPE error takes the form

$$|\hat{v}^\pi - v^\pi| \asymp H^2 \sqrt{\frac{1 + \chi_{\mathcal{Q}}^2(\mu^\pi, \bar{\mu})}{N}} + o(1/\sqrt{N}),$$

where μ^π is some long-term state-action occupancy measure of the target policy π and $\bar{\mu}$ is the data distribution, $\chi_{\mathcal{Q}}^2$ is a variant of χ^2 -divergence *restricted to the family \mathcal{Q}* :

$$\chi_{\mathcal{Q}}^2(p_1, p_2) := \sup_{f \in \mathcal{Q}} \frac{\mathbb{E}_{p_1}[f(x)]^2}{\mathbb{E}_{p_2}[f(x)^2]} - 1.$$

The term $\chi_{\mathcal{Q}}^2(\mu^\pi, \bar{\mu})$ captures the distributional mismatch, between the behavior policy and the target policy, that is relevant to the function class \mathcal{Q} . It determines the theoretical limits of OPE within this function class. In the tabular case, it relates to the worst-case density ratio, which often shows up in importance sampling methods. However, when we use function approximation, this $\chi_{\mathcal{Q}}^2$ divergence term can be significantly smaller than the worst-case density ratio. In particular, our analysis shows that $\chi_{\mathcal{Q}}^2(\mu^\pi, \bar{\mu})$ is the condition number of a finite matrix, which can be reliably estimated. This result suggests that OPE could be more data-efficient with appropriate function approximation.

A summary of technical results of this paper:

- A regression-based algorithm that unifies FQI and plug-in estimation. It does not require knowledge of the behavior policy $\bar{\pi}$, or try to estimate it. It uses iterative regression but does not require Monte Carlo sampling. In the case of linear models, the estimator can be computed easily using simple matrix-vector operations.
- Finite-sample error upper bound for the regression-based policy evaluator. Despite that regression may be biased for OPE, we show that the curse of horizon does not occur as long as $N = \Omega(dH^3)$. A key to the analysis is the use of contraction properties of a Markov process to show that estimation error accumulates linearly in multi-step policy evaluation, instead of exponentially.
- A minimax error lower bound that sets the statistical limit for OPE with function approximation. The lower bound nearly matches our upper bound, therefore proves the efficiency of regression-based FQI.
- A data-dependent confidence bound that can be computed as a byproduct of the FQI algorithm.

1.1. Related Literature

Off-policy policy evaluation (OPE) is often the starting point of batch reinforcement learning. A direct approach is to estimate the transition probability distributions and then execute the target policy on an estimated model. This has been studied in the tabular case with bias and variance analysis (Mannor et al., 2004). In real-world applications, in order to tackle MDPs with infinite or continuous state spaces, one often needs various forms of function approximation, and many methods like fitted Q-iteration and least square policy iteration were developed (Jong & Stone, 2007; Lagoudakis & Parr, 2003; Grunewalder et al., 2012; Fonteneau et al., 2013). Regression methods are often used to fit value functions and to satisfy the Bellman equation (Bertsekas et al., 1995; Sutton & Barto, 2018; Yang & Wang, 2019b).

A popular class of OPE methods use importance sampling (IS) to reweigh sample rewards to get unbiased value estimate of a new policy (Precup, 2000). Doubly robust technique blends IS with model-based estimators to reduce the high variance (Jiang & Li, 2016; Thomas & Brunskill, 2016). Liu et al. (2018) suggested that one should estimate the stationary state occupancy measure instead of the cumulative importance ratio in order to break the curse of horizon. Many IS methods only apply to tabular MDP and require knowledge of the behavior policy. Following these ideas, Nachum et al. (2019) proposed a minimax optimization problem that uses function approximation to learn the IS weights, without requiring knowledge of the behavior policy. Dann et al. (2019) provided error bounds and certificates for the tabular case to achieve accountability. Liu et al. (2019) studied off-policy gradient method for batch data policy optimization.

On the theoretical side, the sharpest OPE error bound to our best knowledge is given by Xie et al. (2019) and Yin & Wang (2020), which applies to time-inhomogeneous, tabular MDP. Jiang & Li (2016) provided a Cramer-Rao lower bound for discrete-tree MDP. To the authors’ best knowledge, most existing theoretical results on OPE apply only to tabular MDP without function approximation. Le et al. (2019) studied batch RL and FQI for policy learning, evaluation and provides generalization bounds that depend on the VC dimension of the function class. Their results require a “concentration coefficient” assumption that the elementwise ratio between generating and target density functions are uniformly bounded across all states, actions and policies. In comparison, our results do not require such concentration condition, and appear to be the first and sharpest error bounds for OPE with linear function approximation.

2. Problem and Model

In this paper, we study off-policy policy evaluation of an Markov decision process (MDP) when we only have a fixed dataset of empirical transitions. An instance of MDP is a controlled random walk over a state space \mathcal{S} , where at each state s , if we pick action $a \in \mathcal{A}$, the system evolves to a random next state s' according to distribution $p(s' | s, a)$ and generates a reward $r' \in [0, 1]$ with $\mathbb{E}[r' | s, a] = r(s, a)$. A policy π specifies a distribution $\pi(\cdot | s)$ for choosing actions conditioned on the current state s .

Our objective is to evaluate the performance of a *target* policy π at a fixed initial distribution ξ_0 , where the transition model p is unknown. The value to be estimated is the expected cumulative reward in an H -horizon episode, given by

$$v^\pi := \mathbb{E}^\pi \left[\sum_{h=0}^H r(s_h, a_h) \mid s_0 \sim \xi_0 \right], \quad (1)$$

where $a_h \sim \pi(\cdot | s_h)$, $s_{h+1} \sim p(\cdot | s_h, a_h)$, \mathbb{E}^π denotes expectation over the sample path generated under policy π .

Let $\mathcal{D} = \{(s_n, a_n, s'_n, r'_n)\}_{n=1}^N$ be a set of sample transitions, where each s'_n is sampled from distribution $p(\cdot | s_n, a_n)$. The sample transitions may be collected from multiple trajectories and under a possibly unknown *behavior* policy denoted as $\bar{\pi}$. Our goal is to estimate v^π from \mathcal{D} .

Given a *target* policy π and a reward function r , the state-action value functions, also known as Q functions, are defined as, for $h = 0, 1, \dots, H$,

$$Q_h^\pi(s, a) := \mathbb{E}^\pi \left[\sum_{h'=h}^H r(s_{h'}, a_{h'}) \mid s_h = s, a_h = a \right], \quad (2)$$

where $a_{h'} \sim \pi(\cdot | s_{h'})$, $s_{h'+1} \sim p(\cdot | s_{h'}, a_{h'})$. Let $\mathcal{X} := \mathcal{S} \times \mathcal{A}$. Define the *conditional transition operator* $\mathcal{P}^\pi : \mathbb{R}^{\mathcal{X}} \rightarrow \mathbb{R}^{\mathcal{X}}$ as

$$\mathcal{P}^\pi f(s, a) := \mathbb{E}^\pi [f(s', a') \mid s, a] \quad \text{for any } f : \mathcal{X} \rightarrow \mathbb{R},$$

where $s' \sim p(\cdot \mid s, a)$ and $a' \sim \pi(\cdot \mid s')$. Throughout the paper, we suppose that \mathcal{P}^π operates in a function class \mathcal{Q} , such that we can approximate unknown Q functions within this family. Assume without loss of generality that $\mathbf{1} \in \mathcal{Q}$.

Assumption 1 (Function class). *For any $f \in \mathcal{Q}$, $\mathcal{P}^\pi f \in \mathcal{Q}$, and $r \in \mathcal{Q}$. It follows that $Q_0^\pi, \dots, Q_H^\pi \in \mathcal{Q}$, where $\mathcal{Q} \subseteq \mathbb{R}^{\mathcal{X}}$.*

In most parts of the paper, we assume that the transition data are collected from multiple independent episodes.

Assumption 2 (Data generating process). *The dataset \mathcal{D} consists of samples from K i.i.d. episodes $\tau_1, \tau_2, \dots, \tau_K$. Each τ_k has H consecutive sample transitions generated by some policy on a single sample path, i.e., $\tau_k = (s_{k,0}, a_{k,0}, r'_{k,0}, s_{k,1}, a_{k,1}, r'_{k,1}, \dots, s_{k,H}, a_{k,H}, r'_{k,H})$.*

We will focus mainly on the case where \mathcal{Q} is a linear space spanned by d feature functions ϕ_1, \dots, ϕ_d . Also note that the behavior policy $\bar{\pi}$ is *not* known.

Notations Denote $\mathcal{X} = \mathcal{S} \times \mathcal{A}$. Let $\mathbb{R}^{\mathcal{X}}$ be the collection of all functions $f : \mathcal{X} \rightarrow \mathbb{R}$. For any $f \in \mathbb{R}^{\mathcal{X}}$, define $f^\pi : \mathcal{S} \rightarrow \mathbb{R}$ by $f^\pi(s) = \int_{\mathcal{A}} f(s, a) \pi(a \mid s) da$. If A is a positive symmetric semidefinite matrix, let $\sigma_{\min}(A)$ denote its smallest eigenvalue, and let $A^{1/2}$ denote the positive symmetric semidefinite matrix that $A^{1/2} A^{1/2} = A$. For nonnegative $\{a_n\}_{n=1}^\infty$ and $\{b_n\}_{n=1}^\infty$, we denote $a_n \lesssim b_n$ if there exists $c > 0$ such that $a_n \leq c b_n$ for $n = 1, 2, \dots$. Let $\{X_n\}_{n=1}^\infty$ be a sequence of random variables and $\{a_n\}_{n=1}^\infty \subseteq \mathbb{R}$ be deterministic. We write $X_n = O_{\mathbb{P}}(a_n)$ if for any $\delta > 0$ there exists $M > 0$ such that $\mathbb{P}(|X_n| > a_n M) \leq \delta$ for all n . If a distribution p is absolutely continuous with respect to distribution q , the Pearson χ^2 -divergence is defined by $\chi^2(p, q) := \mathbb{E}_q \left[\left(\frac{dp}{dq} - 1 \right)^2 \right]$.

3. Regression-Based Off-Policy Evaluation

We consider a fitted Q-iteration method for new policy evaluation using linear function approximation. We show that it is equivalent to a model-based method that estimates a conditional mean operator that embeds the unknown p into the feature space. They admit a simple matrix-vector implementation when \mathcal{Q} is a linear model with finite dimension.

3.1. Fitted Q-iteration (FQI)

The Q-functions satisfy the Bellman equation

$$Q_{h-1}^\pi(s, a) = r(s, a) + \mathbb{E} [V_h^\pi(s') \mid s, a] \quad (3)$$

for $h = 1, 2, \dots, H$, where $s' \sim p(\cdot \mid s, a)$, $V_h^\pi : \mathcal{S} \rightarrow \mathbb{R}$ is the value function defined as $V_h^\pi(s) := \int_{\mathcal{A}} Q_h^\pi(s, a) \pi(a \mid s) da$.

For the given target policy π , we apply regression recursively by letting $Q_{H+1}^\pi := 0$ and for $h = H, H-1, \dots, 0$,

$$\widehat{Q}_h^\pi := \arg \min_{f \in \mathcal{Q}} \left\{ \sum_{n=1}^N \left(f(s_n, a_n) - r'_n - \int_{\mathcal{A}} \widehat{Q}_{h+1}^\pi(s'_n, a) \pi(a \mid s'_n) da \right)^2 + \lambda \rho(f) \right\}, \quad (4)$$

where $\lambda \geq 0$ and $\rho(\cdot)$ is a regularization function. The scheme above provides a recursive way to evaluate $\widehat{Q}_H^\pi, \widehat{Q}_{H-1}^\pi, \dots, \widehat{Q}_0^\pi$ and v^π by regression using empirical data. It is essentially a fitted Q-iteration. The full algorithm is summarized in Algorithm 1.

Algorithm 1 Fitted Q-iteration for Off-Policy Evaluation (FQI-OPE)

Input: initial distribution ξ_0 , target policy π , horizon H , function class \mathcal{Q} , sample transitions $\mathcal{D} = \{(s_n, a_n, s'_n, r'_n)\}_{n=1}^N$

Let $\widehat{Q}_{H+1}^\pi := 0$;

for $h = H, H-1, \dots, 1$ **do**

 Calculate \widehat{Q}_h by solving (4);

end for

Output: $\widehat{v}_{\text{FQI}}^\pi := \int_{\mathcal{X}} \widehat{Q}_0^\pi(s, a) \xi_0(s) \pi(a \mid s) ds da$

3.2. An equivalent model-based method using conditional mean operator

The preceding FQI method can be equivalently viewed as a model-based plug-in estimator. Recall the *conditional transition operator* $\mathcal{P}^\pi : \mathbb{R}^{\mathcal{X}} \rightarrow \mathbb{R}^{\mathcal{X}}$ is

$$\mathcal{P}^\pi f(s, a) := \mathbb{E}^\pi [f(s', a') \mid s, a] \quad \text{for any } f : \mathcal{X} \rightarrow \mathbb{R}.$$

Under Assumption 1, it always holds that $\mathcal{P}^\pi Q_h^\pi \in \mathcal{Q}$. To this end, we are only interested in a ‘‘projection’’ of ground-truth \mathcal{P}^π onto \mathcal{Q} . We estimate the conditional transition operator by $\widehat{\mathcal{P}}^\pi$: for any $f : \mathcal{X} \rightarrow \mathbb{R}$, let

$$\widehat{\mathcal{P}}^\pi f := \arg \min_{g \in \mathcal{Q}} \left\{ \sum_{n=1}^N \left(g(s_n, a_n) - \int_{\mathcal{A}} f(s'_n, a) \pi(a \mid s'_n) da \right)^2 + \lambda \rho(g) \right\}. \quad (5)$$

We can see that, if $N \rightarrow \infty$, $\widehat{\mathcal{P}}^\pi$ converges to a projected version of \mathcal{P}^π onto \mathcal{Q} . Denote $\phi(\cdot) := [\phi_1(\cdot), \dots, \phi_d(\cdot)]^\top : \mathcal{X} \rightarrow \mathbb{R}^d$. In the case where \mathcal{Q} is a linear space given by $\mathcal{Q} = \{\phi(\cdot)^\top w \mid w \in \mathbb{R}^d\}$ and $\rho(\cdot)$ is taken as

$$\rho(f) := \|w\|_2^2 \quad \text{for } f(\cdot) = \phi(\cdot)^\top w, \quad (6)$$

the constructed $\widehat{\mathcal{P}}^\pi$ in (5) corresponds to an estimated \widehat{p} of the form

$$\widehat{p}(\cdot \mid s, a) := \phi(s, a)^\top \widehat{\Sigma}^{-1} \left(\sum_{n=1}^N \phi(s_n, a_n) \delta_{s'_n}(\cdot) \right),$$

where $\widehat{\Sigma} := \lambda I + \sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top$ is the empirical covariance matrix and $\delta_{s'}(\cdot)$ denotes the Dirac measure. Note this \widehat{p} is not necessary a transition kernel.

We adopt a model-based approach and use $\widehat{\mathcal{P}}^\pi$ in the Bellman equation as a plug-in estimator. In particular, let

$$\widehat{r} := \arg \min_{f \in \mathcal{Q}} \left\{ \sum_{n=1}^N (f(s_n, a_n) - r'_n)^2 + \lambda \rho(f) \right\}, \quad (7)$$

and $\widehat{Q}_{H+1}^\pi := 0$,

$$\widehat{Q}_{h-1}^\pi := \widehat{r} + \widehat{\mathcal{P}}^\pi \widehat{Q}_h^\pi, \quad h = H+1, H, \dots, 1.$$

Then we can estimate the policy value by

$$\widehat{v}_{\text{plug-in}}^\pi := \int_{s,a} \widehat{Q}_0^\pi(s, a) \xi_0(s) \pi(a | s) ds da.$$

It is easy to verify that this plug-in estimator is equivalent to the earlier FQI estimator. See the proof in Appendix A.

Theorem 1 (Equivalence between FQI and a model-based method). *If \mathcal{Q} is a linear space and ρ is given by (6), Algorithm 1 and the preceding model-based approach generate identical policy value estimators, i.e. $\widehat{v}^\pi := \widehat{v}_{\text{FQI}}^\pi = \widehat{v}_{\text{plug-in}}^\pi$.*

When \mathcal{Q} is a d -dimensional linear space with feature mapping ϕ , under Assumption 1, there exists a matrix $M^\pi \in \mathbb{R}^{d \times d}$ such that

$$\phi(s, a)^\top M^\pi = \mathbb{E}[\phi^\pi(s')^\top | s, a], \quad \forall (s, a) \in \mathcal{X},$$

where $\phi^\pi(s) := \int \phi(s, a) \pi(a | s) da$. We refer to M^π as the *matrix mean embedding* of the conditional transition operator \mathcal{P}^π . We can implement Algorithm 1 in simple vector forms. We embed the one-step reward function and conditional transition operator into a vector and a matrix, respectively:

$$\widehat{r}(\cdot) = \phi(\cdot)^\top \widehat{R} \quad \text{with } \widehat{R} := \widehat{\Sigma}^{-1} \left(\sum_{n=1}^N r'_n \phi(s_n, a_n) \right),$$

$$\widehat{M}^\pi := \widehat{\Sigma}^{-1} \left(\sum_{n=1}^N \phi(s_n, a_n) \phi^\pi(s'_n)^\top \right).$$

The corresponding conditional mean operator $\widehat{\mathcal{P}}^\pi$ is

$$\widehat{\mathcal{P}}^\pi f(s, a) = \phi(s, a)^\top \widehat{M}^\pi w, \quad \text{for } f(\cdot) = \phi(\cdot)^\top w. \quad (9)$$

We represent \widehat{Q}_h^π in the form of $\widehat{Q}_h^\pi(s, a) = \phi(s, a)^\top \widehat{w}_h^\pi$. In this way, we can easily compute \widehat{Q}_h^π using recursive compact vector-matrix operations, as given in Algorithm 2.

Algorithm 2 Conditional Mean Embedding for Policy Evaluation (CME-PE)

Input: initial distribution ξ_0 , target policy π , horizon H , a basis $\{\phi_1, \dots, \phi_d\}$ of \mathcal{Q} , sample transitions $\mathcal{D} = \{(s_n, a_n, s'_n, r'_n)\}_{n=1}^N$,

Estimate \widehat{R} and \widehat{M}^π according to (8);

Let $\widehat{w}_{H+1}^\pi := 0$;

Let $\nu_0^\pi := \int_{\mathcal{X}} \phi(s, a) \xi_0(s) \pi(a | s) ds da$;

for $h = H, H-1, \dots, 0$ **do**

 Calculate $\widehat{w}_h^\pi := \widehat{R} + \widehat{M}^\pi \widehat{w}_{h+1}^\pi$;

end for

Output: $\widehat{v}^\pi := (\nu_0^\pi)^\top \widehat{w}_0^\pi$

3.3. Relations to other methods

Our method turns out to be closely related to variants of importance sampling method for OPE. For examples:

- *Marginalized importance sampling*: Our FQI estimator takes the form

$$\widehat{v}^\pi = \frac{1}{N} \sum_{n=1}^N \widehat{w}_{\pi/\mathcal{D}}(s_n, a_n) r'_n$$

where

$$\widehat{w}_{\pi/\mathcal{D}}(s, a) := N \sum_{h=0}^H (\nu_0^\pi)^\top (\widehat{M}^\pi)^h \widehat{\Sigma}^{-1} \phi(s, a).$$

By viewing $\widehat{w}_{\pi/\mathcal{D}}(s, a)$ as weights, our estimator can be obtained equivalently by importance sampling. In the special tabular case, our \widehat{v}^π is *equivalent* to the marginalized importance sampling (MIS) estimator in (Yin & Wang, 2020).

- *DualDICE*: Nachum et al. (2019) proposed a minimax formulation to find the stationary state occupancy measure and residue (weight for importance sampling) with function approximation. We observe that, if those function classes are taken to be \mathcal{Q} , a version of DualDICE produces the same estimator as the FQI estimator. The two methods can be viewed as dual to each other.

These intriguing relations permit a unified view of OPE methods. See Appendix A for more discussions.

4. Finite-Sample Error Bound

Recall that \mathcal{D} is a collection of K independent H -horizon trajectories. Let Σ be the uncentered covariance matrix of the data distribution:

$$\Sigma = \mathbb{E} \left[\frac{1}{H} \sum_{h=0}^{H-1} \phi(s_{1,h}, a_{1,h}) \phi(s_{1,h}, a_{1,h})^\top \right],$$

which is determined by the unknown behavior policy $\bar{\pi}$. Given a target policy π , let ξ^π be an invariant distribution of the Markov chain with transition kernel $p^\pi(s' | s) = \int_{\mathcal{A}} p(s' | s, a) \pi(a | s) da$. Define

$$\Sigma^\pi := \mathbb{E}[\phi^\pi(s) \phi^\pi(s)^\top | s \sim \xi^\pi].$$

We assume $\phi(s, a)^\top \Sigma^{-1} \phi(s, a) \leq C_1 d$ without loss of generality. Theorem 2 provides an instance-dependent policy evaluation upper bound. Its complete proof is given in Appendix B.

Theorem 2 (Upper bound). *Let $\delta \in (0, 1)$. Under Assumptions 1 and 2, if $N \geq 20\kappa_1(2 + \kappa_2)^2 \ln(12dH/\delta) C_1 d H^3$ and $\lambda \leq \ln(12dH/\delta) C_1 d H \sigma_{\min}(\Sigma)$, then with probability at least $1 - \delta$,*

$$\begin{aligned}
 & |v^\pi - \hat{v}^\pi| \\
 & \leq \sum_{h=0}^H (H-h+1) \sup_{f \in \mathcal{Q}} \frac{\mathbb{E}^\pi [f(s_h, a_h) \mid s_0 \sim \xi_0]}{\sqrt{\mathbb{E}[\frac{1}{H} \sum_{h=0}^{H-1} f^2(s_{1,h}, a_{1,h})]}} \\
 & \quad \cdot \sqrt{\frac{\ln(12/\delta)}{2N}} + \frac{C \ln(12dH/\delta) dH^{3.5}}{N},
 \end{aligned} \tag{10}$$

where $C := 15\kappa_1 C_1 (3 + \kappa_2) \sqrt{(\nu_0^\pi)^\top \Sigma^{-1} \nu_0^\pi}$,
 $\kappa_1 := \text{cond}(\Sigma^{-1/2} \Sigma^\pi \Sigma^{-1/2})$,
 $\kappa_2 := \|\Sigma^{-1/2} \mathbb{E}[\frac{1}{H} \sum_{h=1}^H \phi^\pi(s_{1,h}) \phi^\pi(s_{1,h})^\top] \Sigma^{-1/2}\|_2 \vee 1$.

Additionally, if either one of the following holds:

- $\phi(s, a)^\top \Sigma^{-1} \phi(s', a') \geq 0$ for any $(s, a), (s', a') \in \mathcal{X}$;
- the MDP is time-inhomogeneous,

the upper bound can be improved to

$$\begin{aligned}
 & |v^\pi - \hat{v}^\pi| \\
 & \leq \sup_{f \in \mathcal{Q}} \frac{\mathbb{E}^\pi [\sum_{h=0}^H (H-h+1) f(s_h, a_h) \mid s_0 \sim \xi_0]}{\sqrt{\mathbb{E}[\frac{1}{H} \sum_{h=0}^{H-1} f^2(s_{1,h}, a_{1,h})]}} \\
 & \quad \cdot \sqrt{\frac{\ln(12/\delta)}{2N}} + \frac{C \ln(12dH/\delta) dH^{3.5}}{N}.
 \end{aligned} \tag{11}$$

Distributional mismatch as a \mathcal{Q} - χ^2 -divergence.

Let $\bar{\mu}$ be the expected occupancy measure of observation $\{(s_n, a_n)\}_{n=1}^N$. Let μ^π be the weighted occupancy distribution of (s_h, a_h) under policy π and ξ_0 , given by

$$\mu^\pi(s, a) := \frac{\mathbb{E}^\pi [\sum_{h=0}^H (H-h+1) \mathbf{1}(s_h = s, a_h = a)]}{\sum_{h=0}^H (H-h+1)}.$$

The upper bound (11) can be simplified to

$$|\hat{v}^\pi - v^\pi| \leq CH^2 \sqrt{\frac{1 + \chi_{\mathcal{Q}}^2(\mu^\pi, \bar{\mu})}{N}} + O(N^{-1}).$$

Moreover, each mismatch term in (10) has a vector form

$$\frac{\mathbb{E}^\pi [f(s_h, a_h) \mid s_0 \sim \xi_0]}{\sqrt{\mathbb{E}[\frac{1}{H} \sum_{h=0}^{H-1} f^2(s_{1,h}, a_{1,h})]}} = \sqrt{(\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi},$$

where $\nu_h^\pi := \mathbb{E}^\pi [\phi(s_h, a_h) \mid s_0 \sim \xi_0]$, so it can be estimated tractably.

The case of tabular MDP.

In the tabular case, the condition $\phi(s, a)^\top \Sigma^{-1} \phi(s', a') \geq 0$ holds for all $(s, a), (s', a') \in \mathcal{X}$. It can be easily seen that the error bound (11) has a strong connection with the χ^2 -divergence between the state-action distributions under the behavior and target policies.

Corollary 1 (Upper bound in tabular case). *In the tabular case with $\mathcal{Q} = \mathbb{R}^{\mathcal{X}}$, if N is sufficiently large and $\lambda = 0$, then with probability at least $1 - \delta$,*

$$|v^\pi - \hat{v}^\pi| \leq 3H^2 \sqrt{1 + \chi^2(\mu^\pi, \bar{\mu})} \sqrt{\frac{\ln(12/\delta)}{2N}} + O(N^{-1}), \tag{12}$$

where $\chi^2(\cdot, \cdot)$ denotes the Pearson χ^2 -divergence. If the MDP is also time-inhomogeneous, then

$$\begin{aligned}
 |v^\pi - \hat{v}^\pi| & \leq \sqrt{H \sum_{h=0}^H \sum_{s,a} \frac{\mu_h^\pi(s, a)^2}{\bar{\mu}_h(s, a)} \text{Var}[r' + V_{h+1}^\pi(s') \mid s, a]} \\
 & \quad \cdot \sqrt{\frac{2 \ln(12/\delta)}{N}} + o(N^{-1/2}),
 \end{aligned} \tag{13}$$

where $\bar{\mu}_h$ is the marginal distribution of $(s_{1,h}, a_{1,h})$ and μ_h^π is the marginal distribution of (s_h, a_h) under policy π and ξ_0 .

The tabular-case upper bound (13) has the same form with Theorem 3.1 in Yin & Wang (2020). The proof of Corollary 1 is deferred to Appendix B.7.

4.1. Proof Outline

We decompose the error into three terms: $v^\pi - \hat{v}^\pi = E_1 + E_2 + E_3$, where E_1 is a linear function of $\hat{\mathcal{P}}^\pi - \mathcal{P}^\pi$, E_2 is a high-order function of $\hat{\mathcal{P}}^\pi - \mathcal{P}^\pi$ and $E_3 = O(\lambda)$. In the following, we outline the analysis of E_1 and E_2 .

First-order term E_1 : This linear error term takes the form $E_1 = \frac{1}{N} \sum_{n=1}^N e_n$, where

$$\begin{aligned}
 e_n & := \sum_{h=0}^H (\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \\
 & \quad \cdot \left(Q_h^\pi(s_n, a_n) - (r'_n + V_{h+1}^\pi(s'_n)) \right).
 \end{aligned}$$

Define a filtration $\{\mathcal{F}_n\}_{n=1, \dots, N}$ where \mathcal{F}_n is generated by $(s_1, a_1, s'_1, r'_1), \dots, (s_{n-1}, a_{n-1}, s'_{n-1}, r'_{n-1})$ and (s_n, a_n) . Then e_1, e_2, \dots, e_N is a martingale difference sequence with respect to $\{\mathcal{F}_n\}_{n=1, \dots, N}$. In what is next, we analyze $\text{Var}[e_n \mid \mathcal{F}_n]$ and apply the Freedman's inequality (Freedman, 1975) to derive a finite sample upper bound for E_1 .

Consider the conditional variance $\text{Var}[e_n \mid \mathcal{F}_n]$. By using the Cauchy-Schwarz inequality and the relation $\text{Var}[r'_n + V_{h+1}^\pi(s'_n) \mid s_n, a_n] \leq \frac{1}{4}(H-h+1)^2$, we have

$$\begin{aligned}
 \text{Var}[e_n \mid \mathcal{F}_n] & = \mathbb{E}[e_n^2 \mid s_n, a_n] \\
 & \leq \frac{1}{4} \left(\sum_{h=0}^H (H-h+1) \sqrt{(\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi} \right) \\
 & \quad \cdot \left(\sum_{h=0}^H \frac{H-h+1}{\sqrt{(\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi}} \left((\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \right)^2 \right).
 \end{aligned} \tag{14}$$

We learn from matrix-form Bernstein inequality that $\frac{1}{N} \sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top$ concentrates around Σ with high probability. It follows that

$$\begin{aligned}
 & \sum_{n=1}^N \left((\nu_h^\pi)^\top \Sigma^{-1} \phi(s_n, a_n) \right)^2 \\
 &= (\nu_h^\pi)^\top \Sigma^{-1} \left(\sum_{n=1}^N \phi(s_n, a_n) \phi(s_n, a_n)^\top \right) \Sigma^{-1} \nu_h^\pi \quad (15) \\
 &= (\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi (N + \sqrt{dH} \cdot O_{\mathbb{P}}(\sqrt{N})).
 \end{aligned}$$

Plugging (15) into (14) and taking the summation, we obtain

$$\begin{aligned}
 \sum_{n=1}^N \text{Var}[e_n | \mathcal{F}_n] &\leq \frac{1}{4} \left(\sum_{h=0}^H (H - h + 1) \sqrt{(\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi} \right)^2 \\
 &\quad \cdot (N + \sqrt{dH} \cdot O_{\mathbb{P}}(\sqrt{N})).
 \end{aligned}$$

It follows from the Freedman's inequality that with high probability,

$$|E_1| \lesssim \frac{1}{\sqrt{N}} \sum_{h=0}^H (H - h + 1) \sqrt{(\nu_h^\pi)^\top \Sigma^{-1} \nu_h^\pi} + \frac{\sqrt{dH}}{N}.$$

High-order term E_2 (bias-inducing term): The high-order term E_2 involves powers of $\widehat{\mathcal{P}}^\pi - \mathcal{P}^\pi$. We use the contraction property of Markov process with respect to its invariant measure, in particular,^f

$$\left\| (\Sigma^\pi)^{1/2} M^\pi (\Sigma^\pi)^{-1/2} \right\|_2 \leq 1. \quad (16)$$

where $\Sigma^\pi = \mathbb{E}[\phi^\pi(s) \phi^\pi(s)^\top | s \sim \xi^\pi]$, ξ^π is an invariant distribution under policy π . Assume Σ^π has full rank for simplicity.

By using the contraction property, we will see that the value error will not grow exponentially in H for large N . We have:

$$\begin{aligned}
 |E_2| &\leq \sum_{h=0}^H \sqrt{(\nu_0)^\top (\Sigma^\pi)^{-1} \nu_0^\pi} \cdot \text{Err}(Q_h^\pi) \\
 &\quad \cdot \left((1 + \text{Err}(\widehat{M}^\pi))^h (1 + \text{Err}(N\widehat{\Sigma}^{-1})) - 1 \right), \quad (17)
 \end{aligned}$$

where the explicit definitions of errors $\text{Err}(\widehat{M}^\pi)$, $\text{Err}(N\widehat{\Sigma}^{-1})$ and $\text{Err}(Q_h^\pi)$ can be found in Lemma B.7, Appendix B.4. By concentration arguments, we can show $\text{Err}(\widehat{M}^\pi)$, $\text{Err}(N\widehat{\Sigma}^{-1}) \lesssim \sqrt{dH/N}$ and $\text{Err}(Q_h^\pi) \lesssim (H - h + 1) \sqrt{d/N}$ with high probability. According to (17), as long as $\text{Err}(\widehat{M}^\pi) \lesssim H^{-1}$, the policy evaluation error will not grow exponentially in H . As a result, if $N \gtrsim dH^3$, we have $|E_2| \lesssim dH^{3.5}/N$. \square

5. Minimax Lower Bound

In this section, we establish a minimax lower bound that characterizes the hardness of off-policy evaluation using linear function approximators. Theorem 3 nearly matches with the finite-sample upper bound given by Theorem 2. The complete proof of Theorem 3 is given in Appendix C.

Theorem 3 (Minimax lower bound). *Suppose that an MDP instance $M = (p, r)$ satisfies:*

- *There exists a set of high-value states $\bar{\mathcal{S}} \subseteq \mathcal{S}$ and a set of low-value states $\underline{\mathcal{S}} \subseteq \mathcal{S}$ under the target policy π such that $V_h^\pi(s) \geq \frac{3}{4}(H - h + 1)$ if $s \in \bar{\mathcal{S}}$ and $V_h^\pi(s) \leq \frac{1}{4}(H - h + 1)$ if $s \in \underline{\mathcal{S}}$;*
- $\bar{p} := \int_{\bar{\mathcal{S}}} \min_{s \in \mathcal{S}} p^\pi(s' | s) ds' \geq c$ and $\underline{p} := \int_{\underline{\mathcal{S}}} \min_{s \in \mathcal{S}} p^\pi(s' | s) ds' \geq c$ for $c > 0$.¹

For any behavior policy $\bar{\pi}$, when N is sufficiently large, one has

$$\begin{aligned}
 & \inf_{\widehat{v}^\pi} \sup_{M' \in \mathcal{N}(M)} \mathbb{P}_{M'} \left(\left| v^\pi - \widehat{v}^\pi(\mathcal{D}) \right| \geq \frac{\sqrt{c}}{24\sqrt{N}} \right. \\
 & \left. \cdot \sup_{f \in \mathcal{Q}} \frac{\mathbb{E}_{M'}^\pi \left[\sum_{h=0}^{H-1} (H - h) f(s_h, a_h) \mid s_0 \sim \xi_0 \right]}{\sqrt{\mathbb{E}_{M'} \left[\frac{1}{H} \sum_{h=0}^{H-1} f^2(s_{1,h}, a_{1,h}) \right]}} \right) \geq \frac{1}{6}, \quad (18)
 \end{aligned}$$

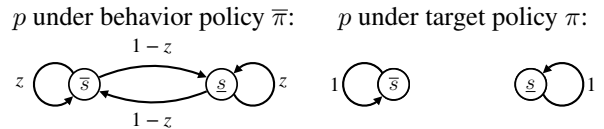
where $\mathcal{N}(M)$ is a small neighborhood of M given by $\mathcal{N}(M) := \{M' = (p', r) \mid \sup_{(s,a) \in \mathcal{X}} \|p'(\cdot | s, a) - p(\cdot | s, a)\|_{\text{TV}} \leq \varepsilon\}$ ($\|\cdot\|_{\text{TV}}$ denotes the total variation, $\varepsilon \gtrsim \sqrt{cd/N}$). $\mathbb{P}_{M'}$ is the probability space of M' , $\widehat{v}^\pi(\mathcal{D})$ is the output of some algorithm \widehat{v}^π when \mathcal{D} is given as the input.

Remark. The minimax lower bound is a worst-case error lower bound that applies to *any estimator*, biased or unbiased. Typical minimax lower bound takes the form of $\inf_{\widehat{v}^\pi} \sup_{\mathcal{M}}$ where the sup is taken over the entire class of MDP \mathcal{M} . Our lower bound is much stronger and can be easily relaxed to the typical form.

Compare Theorems 2 and 3. They nearly match each other, implying that the \mathcal{Q} - χ^2 -divergence term $\chi_{\mathcal{Q}}^2(\bar{\mu}, \mu^\pi)$ determines the statistical complexity of OPE.

An example. Suppose that there is a high-value state \bar{s} and a low-value state \underline{s} , which are two absorbing states under the target policy π , with reward 1 and 0 respectively.

We construct ϕ , π and $\bar{\pi}$ such that $\phi^\pi(\bar{s}) = [z, 1 - z]^\top$, $\phi^\pi(\underline{s}) = [1 - z, z]^\top$; and $\phi^\pi(\bar{s}) = [1, 0]^\top$, $\phi^\pi(\underline{s}) = [0, 1]^\top$. Here $z \in [0, 1]$ is a parameter. We construct the transition model as:



Suppose that the behavior policy $\bar{\pi}$ initiates at either one of the states with probability $1/2$, and the target policy π

¹We assume the behavior policy $\bar{\pi}$ is deterministic only for the sake of notational simplicity.

always initiates at state \bar{s} . We can see that

$$\Sigma = \begin{bmatrix} z^2 - z + \frac{1}{2} & z(1-z) \\ z(1-z) & z^2 - z + \frac{1}{2} \end{bmatrix},$$

and $\nu_0^\pi = \nu_1^\pi = \dots = \nu_{H-1}^\pi = [1, 0]^\top$. For $z \in [\frac{1}{4}, \frac{3}{4}]$, the distributional mismatch term controlling the lower bound becomes

$$\Theta(H^2) \sqrt{1 + \frac{1}{(2z-1)^2}},$$

where z quantifies how much one can tell apart the two states under the target policy π using data generated by $\bar{\pi}$. When $z \approx 1/2$, one can not distinguish \bar{s} and \underline{s} from data generated by $\bar{\pi}$, where the lower bound becomes unbounded.

5.1. Proof Outline

We start with an arbitrary MDP M with transition kernel p that satisfies the assumption. We will construct a perturbed instance $\tilde{p} = p + \Delta p$ so that the two transition models are similar but have a gap in their policy values, denoted by v^π and \tilde{v}^π .

Construct the perturbation Δp such that $\Delta p(s' | s, a) \geq 0$ if $s' \in \bar{\mathcal{S}}$, $\Delta p(s' | s, a) \leq 0$ if $s' \in \underline{\mathcal{S}}$ and $\Delta p(s' | s, a) = 0$ elsewhere. In particular, we construct the perturbation as

$$\begin{aligned} \Delta p(s' | s, a) &= \phi(s, a)^\top \Delta q(s'), \\ \Delta q(s') &:= \mathbf{x} \cdot \min_{s \in \mathcal{S}} p^\pi(s' | s) \cdot (\underline{p}\mathbb{1}_{\bar{\mathcal{S}}}(s') - \bar{p}\mathbb{1}_{\underline{\mathcal{S}}}(s')), \end{aligned} \quad (19)$$

where \bar{p} and \underline{p} are picked such that $\int_{\mathcal{S}} \Delta p(s' | s, a) ds' = 0$ for any s, a , \mathbf{x} is a vector to be picked later.

Reduction to likelihood test We define likelihood functions $\mathcal{L}(\mathcal{D})$ and $\tilde{\mathcal{L}}(\mathcal{D})$ of transition kernels p and \tilde{p} . The likelihood ratio $\frac{\tilde{\mathcal{L}}(\mathcal{D})}{\mathcal{L}(\mathcal{D})} = \prod_{n=1}^N \frac{\tilde{p}(s'_n | s_n, a_n)}{p(s'_n | s_n, a_n)}$ reflects how likely the observation \mathcal{D} comes from model \tilde{p} rather than p . When $p \approx \tilde{p}$, with high probability, the dataset \mathcal{D} generated by model p has a relatively large likelihood ratio, so that it is hard to distinguish p and \tilde{p} based on observation \mathcal{D} . We prove by a martingale concentration argument that, when N is sufficiently large,

$$\ln \left(\frac{\tilde{\mathcal{L}}(\mathcal{D})}{\mathcal{L}(\mathcal{D})} \right) \gtrsim -\sqrt{N} \sqrt{\mathbf{x}^\top \Sigma \mathbf{x}} - N \cdot \mathbf{x}^\top \Sigma \mathbf{x}$$

with high probability. In particular, we have

$$\mathbb{P} \left(\frac{\tilde{\mathcal{L}}(\mathcal{D})}{\mathcal{L}(\mathcal{D})} \geq \frac{1}{2} \right) \geq \frac{1}{2}. \quad (20)$$

when $\sqrt{\mathbf{x}^\top \Sigma \mathbf{x}} \lesssim N^{-1/2}$. If we further have $|v^\pi - \tilde{v}^\pi| \geq \rho + \tilde{\rho}$ for some constant gaps $\rho, \tilde{\rho} \geq 0$, condition (20) implies that for an arbitrary algorithm \hat{v}^π , only one of the following must hold: either $\mathbb{P}(|v^\pi - \hat{v}^\pi(\mathcal{D})| \geq \rho) \geq \frac{1}{6}$ or $\mathbb{P}(|\tilde{v}^\pi - \hat{v}^\pi(\mathcal{D})| \geq \tilde{\rho}) \geq \frac{1}{6}$. In other words, no algorithm can achieve small OPE error for both p and \tilde{p} .

Constructing similar instances with a gap in values

We have

$$\tilde{v}^\pi - v^\pi = \sum_{h=0}^H \xi_0^\top (\tilde{\mathcal{P}}^\pi)^h (\tilde{\mathcal{P}}^\pi - \mathcal{P}^\pi) Q_{h+1}^\pi. \quad (21)$$

By first-order Taylor expansion and our construction, if the perturbation Δp is sufficiently small, we have

$$\begin{aligned} \tilde{v}^\pi - v^\pi &\approx \sum_{h=0}^H \xi_0^\top (\mathcal{P}^\pi)^h (\tilde{\mathcal{P}}^\pi - \mathcal{P}^\pi) Q_{h+1}^\pi \\ &\gtrsim \sum_{h=0}^{H-1} (H-h) (\nu_h^\pi)^\top \mathbf{x}. \end{aligned} \quad (22)$$

For a given N , we maximize the above value over \mathbf{x} under the constraint $\sqrt{\mathbf{x}^\top \Sigma \mathbf{x}} \lesssim N^{-1/2}$. Then we obtain $\mathbf{x}^* = \frac{c_0 \mathbf{x}_0}{\sqrt{N} \sqrt{\mathbf{x}_0^\top \Sigma \mathbf{x}_0}}$ where $c_0 > 0$ is a constant and $\mathbf{x}_0 = \Sigma^{-1} \sum_{h=0}^{H-1} (H-h) \nu_h^\pi$. In this way, we have shown that $\tilde{v}^\pi - v^\pi \gtrsim \frac{1}{\sqrt{N}} \left\| \sum_{h=0}^{H-1} (H-h) \nu_h^\pi \right\|_{\Sigma^{-1}}^2$ using the above construction of \mathbf{x}^* .

Similarly, one can show that for N sufficiently large, $v^\pi - \tilde{v}^\pi \geq \rho + \tilde{\rho}$ for $\rho = \frac{\sqrt{c}}{24\sqrt{N}} \left\| \sum_{h=0}^{H-1} (H-h) \nu_h^\pi \right\|_{\Sigma^{-1}}^2$ and $\tilde{\rho} = \frac{\sqrt{c}}{24\sqrt{N}} \left\| \sum_{h=0}^{H-1} (H-h) \tilde{\nu}_h^\pi \right\|_{\tilde{\Sigma}^{-1}}^2$, where $\tilde{\nu}_h^\pi$ and $\tilde{\Sigma}$ are counterparts of ν_h^π and Σ under the perturbed model \tilde{p} . Finally, we apply the result of the likelihood test and complete the proof. \square

6. A Computable Confidence Bound

Next we study how to quantify the uncertainty in the policy evaluator given by Algorithm 1. In this section, we assume that the dataset is an arbitrary set of experiences, not necessarily independent episodes. We only assume that the transition samples $\mathcal{D} = \{(s_n, a_n, s'_n, r'_n)\}_{n=1, \dots, N}$ are collected in time order.

Assumption 3. *The dataset \mathcal{D} consists of sample transitions $\{(s_t, a_t, s'_t, r'_t)\}_{t=1}^N$ generated in time order, i.e. adapted to a filtration $\{\mathcal{F}_t\}_{t=1}^N$, where $\{(s_\tau, a_\tau, s'_\tau, r'_\tau)\}_{\tau=1}^t$ are \mathcal{F}_t -measurable.*

Assumption 3 is much weaker than Assumption 2. It allows the samples to be generated from a long single path possibly under a nonstationary adaptive policy, as is typical in online reinforcement learning.

Under this mildest assumption, we provide a confidence bound for the policy evaluation error $|v^\pi - \hat{v}^\pi|$, which can be analytically computed from the data \mathcal{D} .

Theorem 4 (Computable confidence bound). *Let Assumptions 1 and 3 hold. Let $\omega := \max \{\|w\|_2 \mid 0 \leq \phi(s, a)^\top w \leq 1, \forall (s, a) \in \mathcal{X}\}^2$. Assume $\|\phi(s, a)\|_2 \leq 1$*

²Such ω always exists and can be computed priorly

for any $(s, a) \in \mathcal{X}$. For a target policy π , with probability at least $1 - \delta$, we have

$$|v^\pi - \hat{v}^\pi| \leq \sum_{h=0}^H (H - h + 1) \sqrt{(\hat{v}_h^\pi)^\top \hat{\Sigma}^{-1} \hat{v}_h^\pi} \cdot \left(\sqrt{2\lambda\omega + 2\sqrt{2d \ln\left(1 + \frac{N}{\lambda d}\right) \ln\left(\frac{3N^2H}{\delta}\right)}} + \frac{4}{3} \ln\left(\frac{3N^2H}{\delta}\right) \right), \quad (23)$$

where \hat{v}_h^π is given by $(\hat{v}_h^\pi)^\top := (\nu_0^\pi)^\top (\widehat{M}^\pi)^h$.

The proof begins with a decomposition of error given by $v^\pi - \hat{v}^\pi = \sum_{h=0}^H (\hat{v}_h^\pi)^\top (w_h^\pi - (\widehat{R} + \widehat{M}^\pi w_{h+1}^\pi))$, from which we derive

$$|v^\pi - \hat{v}^\pi| \leq \sum_{h=0}^H \sqrt{(\hat{v}_h^\pi)^\top \hat{\Sigma}^{-1} \hat{v}_h^\pi} \cdot \left\| \hat{\Sigma}^{1/2} (w_h^\pi - (\widehat{R} + \widehat{M}^\pi w_{h+1}^\pi)) \right\|_2. \quad (24)$$

We analyze the concentration of $\Theta_h := \left\| \hat{\Sigma}^{1/2} (w_h^\pi - (\widehat{R} + \widehat{M}^\pi w_{h+1}^\pi)) \right\|_2^2$ using a martingale argument that is similar to the bandit literature (e.g., proof of Theorem 5 in (Dani et al., 2008)). The complete proof is given in Appendix D.

The confidence bound given in Theorem 4 can be easily calculated as a byproduct of FQI-OPE (Algorithm 1), since $\hat{v}_h^\pi, \hat{\Sigma}$ were already computed in the iterations. In practice, one can tune the value of λ to get the smallest possible confidence bound.

7. Summary

This paper studies the statistical limits of off-policy evaluation using linear function approximation. We establish a minimax error lower bound that depends on a function class-restricted χ^2 -divergence between the data distribution and the target policy's occupancy measure. We prove that a regression-based FQI method, which can be viewed as a plug-in estimator based on a conditional mean embedding of the transition operator, nearly achieves the minimax lower bound. We also provide a computable confidence bound as a byproduct of the algorithm.

Acknowledgments

Mengdi Wang gratefully acknowledges funding from the U.S. National Science Foundation (NSF) grant CMMI-1653435, Air Force Office of Scientific Research (AFOSR) grant FA9550-19-1-020, and C3.ai DTL.

References

Bertsekas, D. P., Bertsekas, D. P., Bertsekas, D. P., and Bertsekas, D. P. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.

Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. 2008.

Dann, C., Li, L., Wei, W., and Brunskill, E. Policy certificates: Towards accountable reinforcement learning. In *International Conference on Machine Learning*, 2019.

Fonteneau, R., Murphy, S. A., Wehenkel, L., and Ernst, D. Batch mode reinforcement learning based on the synthesis of artificial trajectories. *Annals of operations research*, 208(1):383–416, 2013.

Freedman, D. A. On tail probabilities for martingales. *The Annals of Probability*, 3(1):100–118, 1975.

Grunewalder, S., Lever, G., Baldassarre, L., Pontil, M., and Gretton, A. Modelling transition dynamics in mdps with rkhs embeddings. 2012.

Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, 2016.

Jong, N. K. and Stone, P. Model-based function approximation in reinforcement learning. In *Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems*, pp. 1–8, 2007.

Lagoudakis, M. G. and Parr, R. Least-squares policy iteration. *Journal of machine learning research*, 4(Dec): 1107–1149, 2003.

Le, H. M., Voloshin, C., and Yue, Y. Batch policy learning under constraints. *arXiv preprint arXiv:1903.08738*, 2019.

Liu, Q., Li, L., Tang, Z., and Zhou, D. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pp. 5356–5366, 2018.

Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. Off-policy policy gradient with state distribution correction. In *Conference on Uncertainty in Artificial Intelligence*, 2019.

Mannor, S., Simester, D., Sun, P., and Tsitsiklis, J. N. Bias and variance in value function estimation. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 72, 2004.

Nachum, O., Chow, Y., Dai, B., and Li, L. DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections. In *Advances in Neural Information Processing Systems 32*. 2019.

Precup, D. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, pp. 80, 2000.

- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Thomas, P. and Brunskill, E. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning*, 2016.
- Tropp, J. et al. Freedman’s inequality for matrix martingales. *Electronic Communications in Probability*, 16:262–270, 2011.
- Xie, T., Ma, Y., and Wang, Y.-X. Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. In *Advances in Neural Information Processing Systems*, pp. 9665–9675, 2019.
- Yang, L. F. and Wang, M. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. *arXiv preprint arXiv:1905.10389*, 2019a.
- Yang, L. F. and Wang, M. Sample-optimal parametric q-learning using linearly additive features. In *36th International Conference on Machine Learning, ICML 2019*, pp. 12095–12114. International Machine Learning Society (IMLS), 2019b.
- Yin, M. and Wang, Y.-X. Asymptotically efficient off-policy evaluation for tabular reinforcement learning. *arXiv preprint arXiv:2001.10742*, 2020.