

Estimation and Inference in Distributional Reinforcement Learning

Liangyu Zhang* Yang Peng† Jiadong Liang‡ Wenhao Yang§ Zihua Zhang¶

October 2, 2023

Abstract

In this paper, we study distributional reinforcement learning from the perspective of statistical efficiency. We investigate distributional policy evaluation, aiming to estimate the complete distribution of the random return (denoted η^π) attained by a given policy π . We use the certainty-equivalence method to construct our estimator $\hat{\eta}^\pi$, given a generative model is available. We show that in this circumstance we need a dataset of size $\tilde{O}\left(\frac{|S||A|}{\epsilon^{2p}(1-\gamma)^{2p+2}}\right)$ to guarantee a p -Wasserstein metric between $\hat{\eta}^\pi$ and η^π is less than ϵ with high probability. This implies the distributional policy evaluation problem can be solved with sample efficiency. Also, we show that under different mild assumptions a dataset of size $\tilde{O}\left(\frac{|S||A|}{\epsilon^2(1-\gamma)^4}\right)$ suffices to ensure the Kolmogorov metric and total variation metric between $\hat{\eta}^\pi$ and η^π is below ϵ with high probability. Furthermore, we investigate the asymptotic behavior of $\hat{\eta}^\pi$. We demonstrate that the “empirical process” $\sqrt{n}(\hat{\eta}^\pi - \eta^\pi)$ converges weakly to a Gaussian process in the space of bounded functionals on Lipschitz function class $\ell^\infty(\mathcal{F}_{W_1})$, also in the space of bounded functionals on indicator function class $\ell^\infty(\mathcal{F}_{KS})$ and bounded measurable function class $\ell^\infty(\mathcal{F}_{TV})$ when some mild conditions hold. Our findings give rise to a unified approach to statistical inference of a wide class of statistical functionals of η^π .

1 Introduction

Reinforcement learning has achieved remarkable advancements in various fields, including game-playing [32, 42], robotics systems [17], large language models [27, 26], among others. In classical reinforcement learning which relies on the reward hypothesis [36, 37], one evaluates the performance

*Academy for Advanced Interdisciplinary Studies, Peking University; email: zhangliangyu@pku.edu.cn.

†School of Mathematical Sciences, Peking University; email: pengyang@pku.edu.cn.

‡School of Mathematical Sciences, Peking University; email: jdliang@pku.edu.cn.

§Management Science and Engineering, Stanford University; email: yangwh@stanford.edu.

¶School of Mathematical Sciences, Peking University; email: zhzhhang@math.pku.edu.cn.

of a learning agent by the expected returns (*i.e.*, the expected cumulative sum of a received reward). However, in many applications of reinforcement learning, it is not enough to merely consider the expected returns, because other factors such as uncertainty or risks might also be crucial. For example, when we ask a large language model a question we not only expect it to give a useful answer but want to know how reliable the answer is. An investor should consider the risk-return tradeoff when making investment decisions in financial markets, as high expected returns usually mean greater risks [12]. In the area of healthcare, we are not only interested in the expected performance of a dynamic treatment regime but care about its long-tail performance. Otherwise, it would have the potential to cause serious consequences for patients [18].

Distributional reinforcement learning [24, 1] goes beyond the notion of expected returns and proposes to learn the complete distribution of the random returns. Unlike the classical approach, the distributional perspective offers a comprehensive depiction of the inherent uncertainty (known as aleatoric uncertainty) in the performance of learning agents, due to both the stochastic nature of environments and the actions taken by the agents. By employing the distributional perspective, we might obtain a better understanding of the consequences of the agents' behaviors and have a unified approach for dealing with the issues with regard to risk, uncertainty, and robustness [34, 22]. Indeed, Dabney et al. [8] thought that human brains may also rely on a distributional code of future rewards to make decisions.

In the setting of reinforcement learning, we are usually not fully aware of the stochastic environment and must rely on a dataset to evaluate or train a learning agent. This induces another kind of uncertainty that we call statistical uncertainty (also known as epistemic uncertainty), which stems from limited data. In this paper, we seek to simultaneously address the two kinds of uncertainties aforementioned by developing statistical understandings for distributional reinforcement learning. Specifically, we aim to answer the following two fundamental questions: a) Can we learn the full distribution of random returns in a sample-efficient manner? b) Is it possible to perform statistical inferences from the learned return distributions? We give affirmative answers to both questions with the benefit of a statistical analysis of distributional reinforcement learning presented in our paper. We hope our paper can offer new opportunities in uncertainty quantification in reinforcement learning.

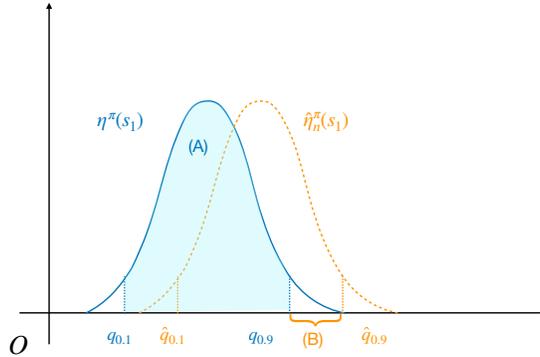


Figure 1. An illustration of two types of uncertainty in RL. Blue distribution: ground-truth return distribution with quantiles $q_{0.1}$ and $q_{0.9}$. Orange distribution: estimated return distribution with quantiles $\hat{q}_{0.1}$ and $\hat{q}_{0.9}$. Shaded area (A): intrinsic uncertainty in RL. Error (B): error caused by statistical uncertainty in RL.

1.1 Our Contributions

In this paper, we focus on the problem of distributional policy evaluation, which lies at the core of distributional reinforcement learning ¹. The goal is to estimate the distribution of random returns of a given policy π written as η^π in an unknown γ -discounted infinite-horizon Markov decision process (MDP). The MDP is assumed to be tabular, *i.e.*, it has finite state space \mathcal{S} and action space \mathcal{A} . Following the common practice in the literature on reinforcement learning, we assume the distribution of the random reward is fully known and the transition probability of that MDP is unknown.

Our estimator $\hat{\eta}^\pi$ is constructed by an empirical version of the distributional dynamic programming. In particular, suppose the underlying MDP is known. Then we may form the distributional Bellman operator \mathcal{T}^π , which is an analog of the classical Bellman operator. Bellemare et al. [1] showed that \mathcal{T}^π is a γ -contraction in the supreme of p -Wasserstein metric with $\eta^\pi := (\eta^\pi(s_1), \dots, \eta^\pi(s_{|\mathcal{S}|}))$ as the unique fixed point. Here $\eta^\pi(s_1)$ means the distribution of returns attained by running policy π given the initial state s_1 . This contracting property leads to the distributional dynamic programming algorithm to compute η^π given that the MDP is already known. However, for the problem of distributional policy evaluation, the underlying MDP is unknown. Therefore, we choose to solve this problem via the certainty-equivalence approach [33, 38]. More specifically, we first build an explicit model of the underlying transition dynamics using a dataset of $n|\mathcal{S}||\mathcal{A}|$ entries obtained by a generative model, and then use that model to construct an empirical version of the

¹Indeed, the control problem in distributional reinforcement learning can be solved by a two-stage procedure. First, we estimate a near-optimal policy $\hat{\pi}$ using some policy learning subroutines. Then it remains to solve a distributional policy evaluation problem, *i.e.* the return distribution of $\hat{\pi}$. See Section 7.3 of [2].

distributional Bellman operator $\widehat{\mathcal{T}}_n^\pi$. An estimator of η^π is then formulated as the fixed point of $\widehat{\mathcal{T}}_n^\pi$, which we denote as $\hat{\eta}_n^\pi := (\hat{\eta}_n^\pi(s_1), \dots, \hat{\eta}_n^\pi(s_{|\mathcal{S}|}))$. Note that we consider a fully non-parametric setting, because $\hat{\eta}_n$ is not restricted in some parametric model and can be any probability distribution.

Our objective is to analyze the statistical performance of the estimated return distribution $\hat{\eta}_n^\pi$. To the best knowledge, we are the first to develop statistical theories for distributional reinforcement learning. Our main contributions are outlined below.

1. We show that under different mild conditions, the distributional dynamic programming algorithm converges to the fixed point when measured by the Kolmogorov-Smirnov metric (KS metric) and total variation metric (TV metric). Interestingly, this convergence occurs despite the distributional Bellman operator no longer being a contraction. Our Findings correct the misconception that distributional dynamic programming is not guaranteed to converge in the KS and TV metric.
2. We provide non-asymptotic bounds for the p -Wasserstein metric, the KS metric and the TV metric between η^π and $\hat{\eta}_n^\pi$. Specifically, we prove $\sup_{s \in \mathcal{S}} W_p(\eta^\pi(s), \hat{\eta}_n^\pi(s)) = \tilde{O}\left(\left[\frac{1}{n(1-\gamma)^{2p+2}}\right]^{1/2p}\right)$, $\sup_{s \in \mathcal{S}} \text{KS}(\eta^\pi(s), \hat{\eta}_n^\pi(s)) = \tilde{O}\left(\frac{1}{n(1-\gamma)^4}\right)$ and $\sup_{s \in \mathcal{S}} \text{TV}(\eta^\pi(s), \hat{\eta}_n^\pi(s)) = \tilde{O}\left(\frac{1}{n(1-\gamma)^4}\right)$ with high probability. Here \tilde{O} means we discard all terms of logarithmic order. Our non-asymptotic results translate to an $\tilde{O}\left(\frac{1}{\epsilon^{2p(1-\gamma)^{2p+2}}}\right)$ complexity bound for the case of the p -Wasserstein metric and $\tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^4}\right)$ complexity bound for the cases of the KS and TV metric, which implies that we can learn the whole distribution of random returns of a given policy in a sample-efficient manner.
3. We give a characterization of the asymptotic behavior of $\hat{\eta}_n^\pi$. We demonstrate that for any $s \in \mathcal{S}$, $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))$ converges weakly to a Gaussian process in $\ell^\infty(\mathcal{F}_{W_1})$. Under different mild conditions, $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))$ also converges weakly to a Gaussian process in $\ell^\infty(\mathcal{F}_{\text{KS}})$ and $\ell^\infty(\mathcal{F}_{\text{TV}})$ for each $s \in \mathcal{S}$. Here \mathcal{F}_{W_1} , \mathcal{F}_{KS} , and \mathcal{F}_{TV} represent the 1-Lipschitz function class, the indicator function class, and the bounded measurable function class respectively. The asymptotic results mentioned earlier enable us to perform statistical inference for η^π . Concretely, we construct asymptotically valid confidence sets for η^π in the forms of W_1 , KS, and TV balls, and asymptotically valid confidence intervals for $\phi(\eta^\pi(s))$, where ϕ can be any Hadamard differentiable functional.
4. At the technical level, our main challenge is that we must work in the infinite-dimensional

space of probability measures. Therefore, most of the techniques developed for classical reinforcement learning theory are not valid anymore. We address the challenge through an analysis of the concentration behaviors as well as asymptotics of $(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)^{-1}(\widehat{\mathcal{T}}_n - \mathcal{T}^\pi)\eta^\pi$. Here $(\mathcal{I} - \mathcal{T}^\pi)^{-1} := \sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$ is defined on a subspace of interest. We achieve this by carefully examining the properties of the distributional Bellman operator \mathcal{T}^π on the vector space of signed measures equipped with different metrics to decouple the dependencies between operators $(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)^{-1}$ and $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)$.

1.2 Related Work

Distributional Reinforcement Learning Distributional reinforcement learning has achieved remarkable success in fields such as communications [14], transportation systems [25], and algorithm discovery [10]. Notable distributional reinforcement learning algorithms include categorical temporal-difference learning [1], quantile temporal-difference learning [7, 6], GAN-based methods [11, 9], actor-critic methods [23], etc. For a comprehensive treatment of distributional reinforcement learning, readers could refer to a very recent book by Bellemare et al. [2].

Despite its empirical success, there is a relative lack of theoretical understanding of distributional reinforcement learning. Rowland et al. [29] analyzed the convergence properties of categorical temporal-difference learning. But their convergence analysis is asymptotic and does not consider sample complexities as well as asymptotic distributions. Recently, Rowland et al. [30] presented similar consistency results for quantile temporal-difference learning. Wu et al. [44] showed that the distribution of returns can be learned by using an algorithm called Fitted Likelihood Estimation given an offline dataset. They also proposed non-asymptotic bounds for the statistical distance between the learned distribution and the ground truth. However, their analysis is modular, assuming that an MLE procedure can achieve good generalization bounds and focusing less on the statistical aspects of distributional policy evaluations. Another line of work treats learning the distribution of returns as an auxiliary task and aims to understand how this auxiliary task can improve policy learning within the framework of classical reinforcement learning. Sun et al. [35] found that such auxiliary tasks can be viewed as a form of regularization and may make the optimization process more stable. Wang et al. [43] explored the statistical benefits of distributional reinforcement learning. They showed that distributional reinforcement learning can yield better non-asymptotic bounds than classical reinforcement learning in the “small loss” scenarios.

Statistical Inference in Reinforcement Learning Statistical inference in the context of reinforcement learning has drawn growing interest in the community. There are a number of works studying the statistical inference problems for expected returns (or value functions). Thomas et al. [39] and Jiang and Li [16] proposed high-confidence bounds for value functions in the setting of off-policy evaluation. Hao et al. [13] devised a bootstrapping procedure to perform statistical inference in off-policy evaluation. Shi et al. [31] modeled the value function with the series/sieve methods and devised confidence intervals for value functions in both the settings of policy evaluation and policy learning. Zhu et al. [46] also constructed asymptotically tight confidence intervals for learned (optimal) value functions. Li et al. [21] and Li et al. [20] considered online statistical inference for value functions in an online reinforcement learning setting.

At the same time, fewer works focus on statistical inferences for other statistical functionals of the return distribution. Yang et al. [45] investigated the asymptotic behaviors of distributionally robust value functions and constructed asymptotically tight confidence bounds. Chandak et al. [5] and Huang et al. [15] proposed methods to estimate the cumulative distribution function (CDF) and confidence band for the ground truth CDF. And statistical inference for statistical functionals can be achieved by the plug-in approach. However, their estimator is based on importance sampling, causing the errors can grow exponentially w.r.t. the horizon length. Also, their confidence intervals are based on non-asymptotic bounds and may thus be too conservative.

The remainder of this paper is organized as follows. In Section 2, we introduce some basic concepts of distributional reinforcement learning. In Section 3, we present our statistical analysis of distributional reinforcement learning. In Section 4, we propose a series of inferential procedures for the return distribution. Section 5 verifies our theoretical findings and tests our inferential procedures through numerical simulations. And Section 6 contains proof outlines of key theoretical results.

2 Distributional Policy Evaluation

Markov Decision Processes An Markov decision process (MDP) is represented by a 5-tuple $M = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}_R, P, \gamma \rangle$, where \mathcal{S} represents a finite state space, \mathcal{A} a finite action space, $\mathcal{P}_R: \mathcal{S} \times \mathcal{A} \rightarrow \Delta([0, 1])$ the distribution of rewards, $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ the transition dynamics, and $\gamma \in (0, 1)$ a discounted factor. Here we use $\Delta(\cdot)$ to represent the set of probability distributions over some set. Given a policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and an initial state $S_0 = s \in \mathcal{S}$, a random trajectory

$\{(S_t, A_t, R_t)_{t=0}^\infty\}$ can be sampled from the MDP using the following procedure:

$$\begin{aligned} A_t | S_t &\sim \pi(\cdot | S_t), \\ R_t | (S_t, A_t) &\sim \mathcal{P}_R(\cdot | S_t, A_t), \\ S_{t+1} | (S_t, A_t) &\sim P(\cdot | S_t, A_t). \end{aligned}$$

We define the return of such trajectories by

$$G^\pi(s) := \sum_{t=0}^{\infty} \gamma^t R_t.$$

Proposition 2.1. *For any policy π and any initial state $s \in \mathcal{S}$, $G^\pi(s): (\Omega, \mathcal{F}, Q) \rightarrow \mathbb{R}$ is a random variable. Here $\Omega = \{\mathcal{S} \times \mathcal{A} \times [0, 1]\}^{\mathbb{N}}$ is the sample space, $\mathcal{F} = \{2^{\mathcal{S}} \times 2^{\mathcal{A}} \times \mathcal{B}[0, 1]\}^{\mathbb{N}}$ is the product σ -field where $\mathcal{B}[0, 1]$ denotes all Borel sets in $[0, 1]$, and Q is a probability measure induced by π, \mathcal{P}_R and P .*

The proof is in Appendix A. Note that we always have $G^\pi(s) \in \left[0, \frac{1}{1-\gamma}\right]$. The expected return $\mathbb{E}G^\pi(s)$ is called the value function and is denoted by $V^\pi(s)$. We also define $\eta^\pi(s)$ as the distribution of $G^\pi(s)$.

Example 2.1 (A Simplest Example of Return Distribution). *Consider a very simple MDP with a single state s_0 and a single action a_0 . We set $\gamma = 1/2$ and $\mathcal{P}_R(\cdot | s_0, a_0) = \text{Bernoulli}(1/2)$. Let π_0 denote the trivial policy $\pi_0(a_0 | s_0) = 1$, then*

$$G^{\pi_0}(s_0) \stackrel{d}{=} \sum_{i=0}^{\infty} \gamma^i X_i, \quad X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(1/2).$$

In other words, we have $\eta^{\pi_0}(s_0) = \text{Uniform}[0, 2]$.

Metrics on the Space of Measures Suppose μ and ν are two probability distributions on \mathbb{R} with finite p -moments ($p \geq 1$). The p -Wasserstein metric between μ and ν is defined as

$$W_p(\mu, \nu) = \left(\inf_{\kappa \in \Gamma(\mu, \nu)} \int_{\mathbb{R}^2} |x - y|^p \gamma(dx, dy) \right)^{1/p}.$$

Elements $\kappa \in \Gamma(\mu, \nu)$ are called couplings of μ and ν , *i.e.*, joint distributions on \mathbb{R}^2 with prescribed marginals μ and ν on each “axis”. Suppose μ and ν have cumulative distribution function F_μ and

F_ν , in the case of $p = 1$ we have

$$W_1(\mu, \nu) = \int_{\mathbb{R}} |F_\mu(x) - F_\nu(x)| dx$$

The Kolmogorov–Smirnov metric (KS metric) is defined as

$$\text{KS}(\mu, \nu) = \sup_{t \in \mathbb{R}} |\mu((-\infty, t]) - \nu((-\infty, t])|.$$

We may bound $\text{KS}(\mu, \nu)$ with $W_1(\mu, \nu)$ when either of μ, ν has bounded densities.

Proposition 2.2. [28, Proposition 1.2] *Assume that $\mu \in \Delta(\mathbb{R})$ has finite moment and μ has a Lebesgue density p_μ that is bounded by C . Then for any $\nu \in \Delta(\mathbb{R})$ with finite moment, $\text{KS}(\mu, \nu) \leq \sqrt{2CW_1(\mu, \nu)}$.*

The total variation distance (TV distance) is defined as

$$\text{TV}(\mu, \nu) = \sup_{A \in \mathcal{B}(\mathbb{R})} |\mu(A) - \nu(A)|,$$

where $\mathcal{B}(\mathbb{R})$ denotes all Borel sets in \mathbb{R} . $\text{TV}(\mu, \nu)$ can also be bounded by $W_1(\mu, \nu)$ when μ and ν have smooth densities.

Proposition 2.3. [4, Theorem 2.1] *Assume that $\mu, \nu \in \Delta(\mathbb{R})$ have Lebesgue densities $p_\mu, p_\nu \in H_1^1(\mathbb{R})$. Specifically, $H_1^1(\mathbb{R})$ represents the L^1 Sobolev space of order 1 defined as*

$$H_1^1(\mathbb{R}) := \{f \in L^1(\mathbb{R}) : D^1 f \in L^1(\mathbb{R}); \|f\|_{H_1^1} = \|f\|_1 + \|D^1 f\|_1 < \infty\}.$$

Here $L^1(\mathbb{R})$ is the space of Lebesgue integrable functions and $\|\cdot\|_1$ is the associated L^1 norm, $D^1 f$ represents the weak derivative of f . Then we have

$$\text{TV}(\mu, \nu) \leq \sqrt{K \left(\|p_\mu\|_{H_1^1} + \|p_\nu\|_{H_1^1} \right) W_1(\mu, \nu)}.$$

Here K is a universal constant.

The 1-Wasserstein metric, the KS metric, and the TV distance are all special cases of integral probability metrics. Specifically, we define

$$\|\mu - \nu\|_{\mathcal{H}} = \sup_{h \in \mathcal{H}} |\mu h - \nu h|,$$

where \mathcal{H} denotes some function class and μh represents $\mathbb{E}_{X \sim \mu} [h(X)]$. If we choose

- $\mathcal{F}_{W_1} := \{f \mid f \text{ is 1-Lipschitz}\}$, then $\|\mu - \nu\|_{\mathcal{F}_{W_1}} = W_1(\mu, \nu)$.
- $\mathcal{F}_{\text{KS}} := \{\mathbb{1}\{\cdot \leq z\} \mid z \in \mathbb{R}\}$, then $\|\mu - \nu\|_{\mathcal{F}_{\text{KS}}} = \text{KS}(\mu, \nu)$.
- $\mathcal{F}_{\text{TV}} := \{f \mid f \text{ is measurable and } \|f\|_\infty \leq 1\}$, then $\|\mu - \nu\|_{\mathcal{F}_{\text{TV}}} = \text{TV}(\mu, \nu)$.

Distributional Bellman Operator It is well-known that the expected returns (also called the value functions) satisfy the Bellman equation. In particular, letting V^π denote $(V^\pi(s_1), \dots, V^\pi(s_{|\mathcal{S}|}))$, we have for any $s \in \mathcal{S}$,

$$\begin{aligned} V^\pi(s) &= [T^\pi(V^\pi)](s) \\ &:= \mathbb{E}_{A \sim \pi(\cdot|s), R \sim \mathcal{P}(\cdot|s, A)} R + \mathbb{E}_{A \sim \pi(\cdot|s), S' \sim P(\cdot|s, A)} V^\pi(S') \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) \int_0^1 r \mathcal{P}_R(dr|s, a) + \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \pi(a|s) P(s'|s, a) V^\pi(s'). \end{aligned} \quad (1)$$

We call the operator $T^\pi: \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$ the Bellman operator. And the Bellman equation suggests that the value function V^π is a fixed point of T^π .

The distributional Bellman equation describes a similar relationship to Equation (1) for the distributions of returns. Letting η^π denote $(\eta^\pi(s_1), \dots, \eta^\pi(s_{|\mathcal{S}|}))$, we have for any $s \in \mathcal{S}$

$$\begin{aligned} \eta^\pi(s) &= [\mathcal{T}^\pi(\eta^\pi)](s) \\ &:= \mathbb{E}_{A \sim \pi(\cdot|s), R \sim \mathcal{P}_R(\cdot|s, A), S' \sim P(\cdot|s, A)} (b_{R, \gamma})_{\#} \eta^\pi(S') \\ &= \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \pi(a|s) P(s'|s, a) \int_0^1 (b_{r, \gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr|s, a). \end{aligned} \quad (2)$$

Here $b_{r, \gamma}: \mathbb{R} \rightarrow \mathbb{R}$ is an affine function defined by $b_{r, \gamma}(x) = r + \gamma x$, and $g_{\#} \mu$ is the push forward measure of μ through function g so that $g_{\#} \mu(A) = \mu(g^{-1}(A))$. The integral $\int_0^1 (b_{r, \gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr|s, a)$ is defined by

$$\left[\int_0^1 (b_{r, \gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr|s, a) \right] (B) = \int_0^1 \left[(b_{r, \gamma})_{\#} \eta^\pi(s') \right] (B) \mathcal{P}_R(dr|s, a)$$

for any Borel set B in $\left[0, \frac{1}{1-\gamma}\right]$. We call the operator $\mathcal{T}^\pi: \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)^{\mathcal{S}} \rightarrow \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)^{\mathcal{S}}$ the distributional Bellman operator, and the distribution of returns η^π is a fixed point of \mathcal{T}^π .

Distributional Dynamic Programming Suppose the MDP M is already known. Given a policy π we may compute the value function V^π by the dynamic programming algorithm. Specifically, assuming $V_{k+1} = T^\pi(V_k)$, we have T^π is a γ -contraction w.r.t. norm $\|\cdot\|_\infty$ on $\mathbb{R}^{\mathcal{S}}$, and thus $\lim_{k \rightarrow \infty} \|V_k - V^\pi\|_\infty = 0$. In analogy to dynamic programming, we may also define distributional dynamic programming, *i.e.*, $\eta^{(k+1)} = \mathcal{T}^\pi \eta^{(k)}$. It can be shown that \mathcal{T}^π is a γ -contraction in the supreme p -Wasserstein metrics. Thus distributional dynamic programming exhibits geometric convergence in the supreme p -Wasserstein metric.

Proposition 2.4. [2, Propositions 15 and 16] *The distributional Bellman operator is a γ -contraction on $\Delta(\mathbb{R})^{\mathcal{S}}$ in the supreme p -Wasserstein metrics. More precisely, for $\eta, \eta' \in \Delta(\mathbb{R})^{\mathcal{S}}$, we have*

$$\sup_{s \in \mathcal{S}} W_p([\mathcal{T}^\pi \eta](s), [\mathcal{T}^\pi \eta'](s)) \leq \gamma \sup_{s \in \mathcal{S}} W_p(\eta(s), \eta'(s)).$$

Furthermore, we have

$$\sup_{s \in \mathcal{S}} W_p(\eta^{(k)}(s), \eta^\pi(s)) \leq \gamma^k \sup_{s \in \mathcal{S}} W_p(\eta^{(0)}(s), \eta^\pi(s))$$

and

$$\lim_{k \rightarrow \infty} \sup_{s \in \mathcal{S}} W_p(\eta^{(k)}(s), \eta^\pi(s)) = 0.$$

When measured by other commonly used probability metrics like the supreme KS metric or the supreme TV distance, the distributional Bellman operator may not be a contraction and distributional dynamic programming may not converge at all [2]. This is because, unlike the cases of dynamic programming, now we operate in an infinite-dimensional space and the metrics may not be equivalent. However, we find that under mild conditions distributional dynamic programming does converge in the supreme KS metric and the supreme TV distance, and the convergences are also geometrically fast. To the best of our knowledge, we are the first to examine the convergence property of distributional dynamic programming w.r.t. the supreme KS metric and the supreme TV distance.

Assumption 1. *Assume that for any $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\mathcal{P}_R(dr | s, a)$ has a Lebesgue density $p_{s,a}^R$ upper-bounded by a constant C .*

Proposition 2.5. *The distributional Bellman operator is non-expansive on $\Delta(\mathbb{R})^{\mathcal{S}}$ in the supreme*

TV metrics. If Assumption 1 holds, then we have

$$\sup_{s \in \mathcal{S}} \text{KS}(\eta^{(k)}(s), \eta^\pi(s)) \leq (\sqrt{\gamma})^k \sup_{s \in \mathcal{S}} \sqrt{CW_1(\eta^{(0)}(s), \eta^\pi(s))}$$

and

$$\lim_{k \rightarrow \infty} \sup_{s \in \mathcal{S}} \text{KS}(\eta^{(k)}(s), \eta^\pi(s)) = 0.$$

To prove Proposition 2.5, we first show that when Assumption 1 is true the distribution of return $\eta^\pi(s)$ must have a bounded density. Then by Proposition 2.2 the KS metric can be controlled by the 1-Wasserstein metric. The full proof can be found in Appendix A.

Assumption 2. Assume that for any $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\mathcal{P}_R(dr \mid s, a)$ has a Lebesgue density $p_{s,a}^R \in H_1^1(\mathbb{R})$ and $\|p_{s,a}^R\|_{H_1^1(\mathbb{R})} \leq M$.

Assumption 2 is indeed strictly stronger than Assumption 1. Precisely, directly applying Sobolev's inequality we have

Proposition 2.6. When Assumption 2 holds, for any $s \in \mathcal{S}$, $a \in \mathcal{A}$, $\sup_{x \in [0, 1/(1-\gamma)]} p_{s,a}^R(x) \leq 2M$.

Proposition 2.7. The distributional Bellman operator is non-expansive on $\Delta(\mathbb{R})^{\mathcal{S}}$ in the supreme TV metrics. If Assumption 2 holds, then we have

$$\sup_{s \in \mathcal{S}} \text{TV}(\eta^{(k)}(s), \eta^\pi(s)) \leq (\sqrt{\gamma})^k \sup_{s \in \mathcal{S}} \sqrt{2MKW_1(\eta^{(0)}(s), \eta^\pi(s))}.$$

and

$$\lim_{k \rightarrow \infty} \sup_{s \in \mathcal{S}} \text{TV}(\eta^{(k)}(s), \eta^\pi(s)) = 0.$$

The proof strategy is similar to that of Proposition 2.5. We first show that when Assumption 2 holds, the $\eta^\pi(s)$ and $\eta^{(k)}(s)$ both have densities in $H_1^1(\mathbb{R})$. Then we can bound the TV distance with the 1-Wasserstein metric via Proposition 2.3. One may refer to Appendix A for the detailed proof.

Empirical Distributional Dynamic Programming We are interested in learning the distribution of returns η^π when the underlying transition dynamic P is unknown and must be estimated from a dataset, which is called the distributional policy evaluation problem. We assume the dataset is generated by a generative model, which is able to return a value of the next state s' following $P(\cdot \mid s, a)$ for any given pair $(s, a) \in \mathcal{S} \times \mathcal{A}$. For each pair $(s, a) \in \mathcal{S} \times \mathcal{A}$, we make n calls to the

generative model and get an array $X_1^{(s,a)}, \dots, X_n^{(s,a)} \stackrel{\text{iid}}{\sim} P(\cdot | s, a)$. Then we may obtain the estimate of the transition probability as

$$\widehat{P}(s' | s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{X_i^{(s,a)} = s'\}.$$

Thus, \widehat{P} defines an empirical MDP $\widehat{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}_R, \widehat{P}, \gamma \rangle$ with the corresponding distribution of returns $\widehat{\eta}_n^\pi$. We also have the following distributional Bellman equation

$$\begin{aligned} \widehat{\eta}_n^\pi(s) &= \left[\widehat{\mathcal{T}}_n^\pi(\widehat{\eta}_n^\pi) \right](s) \\ &:= \mathbb{E}_{A \sim \pi(\cdot | s), R \sim \mathcal{P}(\cdot | s, A), S' \sim \widehat{P}(\cdot | s, A)} (b_{R, \gamma})_{\#} \widehat{\eta}_n^\pi(S') \\ &= \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \pi(a | s) \widehat{P}(s' | s, a) \int_0^1 (b_{r, \gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a), \end{aligned} \tag{3}$$

where $\widehat{\mathcal{T}}_n^\pi$ is called the empirical Bellman operator. Note that $\widehat{\eta}_n^\pi$ may serve as an estimator of the return of distribution η^π and can be computed via the empirical version of distributional dynamic programming, *i.e.*, $\widehat{\eta}^{(k+1)} = \widehat{\mathcal{T}}_n^\pi(\widehat{\eta}^{(k)})$. We devote the remaining parts of the paper to discussions of the statistical properties of $\widehat{\eta}_n^\pi$.

3 Statistical Analysis

In this section, we analyze distributional reinforcement learning from both the non-asymptotic and asymptotic viewpoints. We give the non-asymptotic convergence rates of $\sup_{s \in \mathcal{S}} W_1(\widehat{\eta}_n^\pi(s), \eta^\pi(s))$, $\sup_{s \in \mathcal{S}} \text{KS}(\widehat{\eta}_n^\pi(s), \eta^\pi(s))$, and $\sup_{s \in \mathcal{S}} \text{TV}(\widehat{\eta}_n^\pi(s), \eta^\pi(s))$, which suggest distributional policy evaluation is sample-efficient when a generative model is available. We also study the asymptotics of $\sqrt{n}(\widehat{\eta}_n^\pi(s) - \eta^\pi(s))$ for any $s \in \mathcal{S}$. Under mild conditions, we demonstrate that $\sqrt{n}(\widehat{\eta}_n^\pi(s) - \eta^\pi(s))$ converges weakly to a Gaussian random element in spaces $\ell^\infty(\mathcal{F}_{W_1})$, $\ell^\infty(\mathcal{F}_{\text{KS}})$ and $\ell^\infty(\mathcal{F}_{\text{TV}})$.

3.1 Results on Non-asymptotic Analysis

Our main results of non-asymptotic analysis is given in the following theorems.

Theorem 3.1. *For any fixed policy π ,*

$$\mathbb{E} \sup_{s \in \mathcal{S}} W_1(\widehat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \sqrt{\frac{9 \log |\mathcal{S}|}{n(1 - \gamma)^4}}.$$

And for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2}}{\sqrt{n(1-\gamma)^4}}.$$

To sum up, we show that $n = \tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^4}\right)$ suffices to ensure both $\mathbb{E} \sup_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \epsilon$ and $\sup_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \epsilon$ with high probability, which implies model-based distributional policy evaluation is sample-efficient. The key idea of our proof is that we first analyze the concentration behaviors of the infinite-dimensional operator $(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)$. Then we examine the properties of \mathcal{T}^π in the vector space of signed measures equipped with the W_1 -metric and give a reasonable definition of $(\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1} := \sum_{i=0}^{\infty} (\hat{\mathcal{T}}_n^\pi)^i$ on a product of vector space consisting of signed measures μ such that $\mu\left(\left[0, \frac{1}{1-\gamma}\right]\right) = 0$. This allows us to write $\hat{\eta}_n - \eta^\pi = (\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1} (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$. We can draw the conclusion noting that the operator norm of $(\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1}$ w.r.t. the W_1 -metric is always bounded by $\frac{1}{1-\gamma}$. A more detailed outline of proof can be found in Section 6.

Compared with the minimax optimal $\tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^3}\right)$ sample complexity bound for the model-based policy evaluation [19], our sample complexity bound has an additional $\sqrt{\frac{1}{1-\gamma}}$ factor. In fact, the problem of learning the distribution of returns is harder than the policy evaluation problem because one always has $|V(s) - \hat{V}(s)| \leq W_1(\eta^\pi(s), \hat{\eta}_n^\pi(s))$. However, we speculate the additional $\sqrt{\frac{1}{1-\gamma}}$ factor can be eliminated with more refined analysis techniques specially developed for handling infinite-dimensional cases.

Combine Theorem 3.1 with the elementary inequality $[W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s))]^p \leq \frac{1}{(1-\gamma)^{p-1}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s))$, we can derive the non-asymptotic results for W_p metric.

Corollary 3.1. *For any fixed policy π , $p > 1$,*

$$\mathbb{E} \sup_{s \in \mathcal{S}} W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \left[\frac{9 \log |\mathcal{S}|}{n(1-\gamma)^{2p+2}} \right]^{\frac{1}{2p}}.$$

And for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{s \in \mathcal{S}} W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \left[\frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2}}{\sqrt{n(1-\gamma)^{2p+2}}} \right]^{\frac{1}{p}}.$$

We comment that the slow rate $n^{-\frac{1}{2p}}$ for the W_p -metric is inevitable without assuming additional regularity conditions. Consider an MDP with $\mathcal{S} = \{s_1, s_2, s_3\}$ and $\mathcal{A} = \{a_1\}$. As there is only one

single action a_1 available, the action variable can be safely omitted. We start from s_1 with $r(s_1) = 0$ and $P(s_2 | s_1) = P(s_3 | s_1) = \frac{1}{2}$. And s_2, s_3 are absorbing states with $r(s_2) = 1$ and $r(s_3) = 0$. Suppose after n calls of the generative model we have gained an estimator of $P(s_2 | s_1)$ that is denoted as \hat{p} , then $\eta(s_1) = \frac{1}{2}\delta_0 + \frac{1}{2}\delta_{\frac{1}{1-\gamma}}$, $\hat{\eta}_n(s_1) = (1 - \hat{p})\delta_0 + \hat{p}\delta_{\frac{1}{1-\gamma}}$. We have

$$W_p(\eta(s_1), \hat{\eta}_n(s_1)) = \frac{1}{1-\gamma} \left| \hat{p} - \frac{1}{2} \right|^{\frac{1}{p}}.$$

Since $\hat{p} \sim \text{Binomial}(n, \frac{1}{2})$, $|\hat{p} - \frac{1}{2}|$ is of the order $n^{-\frac{1}{2}}$ by CLT. Thus $W_p(\eta(s_1), \hat{\eta}_n(s_1))$ is of the order $n^{-\frac{1}{2p}}$.

Under Assumption 1, we also have the following bounds on the KS metric.

Theorem 3.2. *Suppose Assumption 1 holds true. For any fixed policy π ,*

$$\mathbb{E} \sup_{s \in \mathcal{S}} \text{KS}(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq C'' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}}.$$

And for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{s \in \mathcal{S}} \text{KS}(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \frac{C' \left(\sqrt{\log |\mathcal{S}|} + \sqrt{\log(1/\delta)} \right)}{n(1-\gamma)^4}.$$

Here C' and C'' are constants only depending on C in Assumption 1.

The upper bound of the supreme KS metric between $\hat{\eta}_n$ and η^π is of the same order as that of the W_1 -metric, which indicates that under mild conditions learning a near-optimal return distribution in the sense of KS metric is not more difficult than learning a near-optimal return distribution in the sense of W_1 -metric. This is somewhat a surprise since the distributional Bellman operator exhibits benign behaviors only when measured by Wasserstein metrics, and for μ, ν with bounded support $W_1(\mu, \nu)$ can be always bounded by $\text{KS}(\mu, \nu)$ multiplying a constant factor.

Simply combining the results in Theorem 3.1 and Proposition 2.2 can only yield a sub-optimal $n^{-\frac{1}{4}}$ convergence rate for $\sup_{s \in \mathcal{S}} \text{KS}(\hat{\eta}_n(s), \eta^\pi(s))$. Instead, we obtain a $n^{-\frac{1}{2}}$ rate using a quite different proof strategy with that of Theorem 3.1. The first challenge is that the operator $(\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1}$ may be unbounded on its domain measured with the norm induced by the KS distance. Therefore, although we can still write $(\hat{\eta}_n^\pi - \eta^\pi) = (\mathcal{T} - \hat{\mathcal{T}}_n^\pi)^{-1} (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$, bounds of $(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$ does not directly translate to bounds of $(\hat{\eta}_n^\pi - \eta^\pi)$. We handle this challenge by an ‘‘expansion trick’’, which

raises yet another technical challenge: we need a stronger notion of concentration of $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$. Specifically, unlike in the proof of Theorem 3.1 where it suffices to bound $W_1(\widehat{\mathcal{T}}_n^\pi \eta^\pi, \mathcal{T}^\pi \eta^\pi)$, here we further need to bound $\text{TV}(\widehat{\mathcal{T}}_n^\pi \eta^\pi, \mathcal{T}^\pi \eta^\pi)$. And we achieve it with an analysis through the lens of density functions. A more detailed proof outline can be found in Section 6.

Theorem 3.3. *Suppose Assumption 2 holds true. For any fixed policy π ,*

$$\mathbb{E} \sup_{s \in \mathcal{S}} \text{TV}(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq K'' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}}.$$

And for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sup_{s \in \mathcal{S}} \text{TV}(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \frac{K' \left(\sqrt{\log |\mathcal{S}|} + \sqrt{\log(1/\delta)} \right)}{\sqrt{n(1-\gamma)^4}}.$$

K', K'' are absolute constants depending only on M in Assumption 2.

Based on the theorem mentioned above, we observe that our upper bounds for the supremum TV distance are of the same order as those for the supremum W_1 metric or the KS distance. Directly applying Theorem 3.1 and Proposition 2.3 only attain a slow $n^{-\frac{1}{4}}$ rate. We rather employ a similar analytical approach as in the case of the KS distance to establish a standard convergence rate of $n^{-\frac{1}{2}}$. We present a more detailed outline of proof in Section 6.

3.2 Results on Asymptotic Analysis

We first give our main results of asymptotic analysis in Theorem 3.4.

Theorem 3.4. *For any fixed policy π , we have for any $s \in \mathcal{S}$*

- (a) $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))$ converge weakly to the process $f \mapsto [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi](s)f$ in $\ell^\infty(\mathcal{F}_{W_1})$, where $\mathcal{F}_{W_1} := \left\{ f \mid f \text{ is supported on } \left[0, \frac{1}{1-\gamma}\right] \text{ and } 1\text{-Lipschitz} \right\}$.
- (b) If Assumption 1 is true, then $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))$ converge weakly to the process $f \mapsto [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi](s)f$ in $\ell^\infty(\mathcal{F}_{\text{KS}})$, where $\mathcal{F}_{\text{KS}} := \left\{ \mathbf{1}_{(-\infty, z]} \mid z \in \left[0, \frac{1}{1-\gamma}\right] \right\}$.
- (c) If Assumption 2 is true, then $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))$ converge weakly to the process $f \mapsto [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi](s)f$ in $\ell^\infty(\mathcal{F}_{\text{TV}})$, where $\mathcal{F}_{\text{TV}} := \left\{ \mathbf{1}_A \mid A \subseteq \left[0, \frac{1}{1-\gamma}\right] \text{ is Borel} \right\}$.

Here the random element $\tilde{\mathbb{G}}^\pi$ is defined as

$$\tilde{\mathbb{G}}^\pi(s) := \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} Z_{s,a,s'} \int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \quad \forall s \in \mathcal{S},$$

where $(Z_{s,a,s'})_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}}$ are zero-mean Gaussians with

$$\text{Cov}(Z_{s_1,a_1,s'_1}, Z_{s_2,a_2,s'_2}) = \mathbb{1}\{(s_1, a_1) = (s_2, a_2)\} P(s'_1 | s_1, a_1) (\mathbb{1}\{s'_1 = s'_2\} - P(s'_2 | s_1, a_1)).$$

And the operator $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$ is defined as $(\mathcal{I} - \mathcal{T}^\pi)^{-1} := \sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$.

At a high level, we depict the asymptotic behavior of $\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi)$ by showing that the “empirical processes” induced by $\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi)$ converge to a Gaussian random element. Moreover, the limiting random element has a simple structure in the sense that it is a linear transformation of a finite mixture of probability distributions with Gaussian coefficients. Our asymptotic results are general in the sense that the conclusions are valid in different spaces under different regularity conditions: $\ell^\infty(\mathcal{F}_{W_1})$, $\ell^\infty(\mathcal{F}_{\text{KS}})$, and $\ell^\infty(\mathcal{F}_{\text{TV}})$. Therefore, our findings have the potential to yield numerous valuable inferential procedures for the field of distributional reinforcement learning. Our proof of Theorem 3.4 builds on the foundation of our non-asymptotic analysis and can be found in Section 3.

4 Statistical Inference

In this section, we consider the statistical inference of distributional reinforcement learning. First, we present non-parametric confidence sets for $\eta^\pi(s)$ in the forms of W_1 , KS, and TV balls. Second, we study inference on Hadamard differentiable functionals, with moments, quantiles, and uniform advantage of policy as special examples.

4.1 Inferences for $\eta^\pi(s)$

Our theoretical findings in Theorems 3.4 allow us to construct confidence sets in the space $\Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)$ for the true return distribution $\eta^\pi(s)$, given any initial state $s \in \mathcal{S}$. Specifically, we can construct three types of confidence sets for $\eta^\pi(s)$: W_1 , KS and TV balls.

Proposition 4.1. *For some fixed policy π and initial state $s \in \mathcal{S}$, define $\rho_1(\alpha) := \frac{z_{1-\alpha/2}}{\sqrt{n}}$, where*

$z_1(p)$ is defined as the p -quantile of $\sup_{f \in \mathcal{F}_{W_1}} \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) f$, Then we have

$$\lim_{n \rightarrow \infty} \mathbb{P} (W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \rho_1(\alpha)) = 1 - \alpha$$

Furthermore, if Assumption 1 holds, we have

$$\sup_{f \in \mathcal{F}_{W_1}} \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) f = \int_0^{\frac{1}{1-\gamma}} \left| \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \mathbb{1} \{ \cdot \leq x \} \right| dx.$$

The proof is in Appendix C. Proposition 4.1 depicts the quantile of $W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s))$. Recall that for two probability distributions μ_1, μ_2 supported on $\left[0, \frac{1}{1-\gamma}\right]$ we have

$$W_1(\mu_1, \mu_2) = \sup_{f \in \mathcal{F}_{W_1}} |\mu_1 f - \mu_2 f| = \int_0^{\frac{1}{1-\gamma}} |F_1(x) - F_2(x)| dx,$$

where F_1 and F_2 are the cumulative distribution functions of μ_1 and μ_2 , respectively. Hence, the asymptotic distribution of $W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s))$ can be described in two different ways using asymptotic results in Theorem 3.4 and the continuous mapping theorem. And $\rho_1(\alpha)$ can be determined accordingly.

Now we are ready to present our confidence sets for $\eta^\pi(s)$ that take the form of 1-Wasserstein balls.

Theorem 4.1. *For some fixed policy π and initial state $s \in \mathcal{S}$, define*

$$C_1(\alpha) := \left\{ \eta \in \Delta \left(\left[0, \frac{1}{1-\gamma}\right] \right) \mid W_1(\eta, \hat{\eta}_n^\pi(s)) \leq \rho_1(\alpha) \right\}.$$

Then we have $\lim_{n \rightarrow \infty} \mathbb{P}(\eta^\pi(s) \in C_1(\alpha)) = 1 - \alpha$.

$C_1(\alpha)$ is asymptotically valid, but it relies on the quantile, *i.e.*, $z_1(1 - \alpha)$, of the unknown limiting distributions that depend on \mathcal{T}^π and η^π . We may get a consistent estimate of $z_1(1 - \alpha)$ by the plug-in approach.

Proposition 4.2. *For any fixed policy π and initial state $s \in \mathcal{S}$, define*

$$\hat{\mathbb{G}}^\pi(s) := \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} \hat{Z}_{s,a,s'} \int_0^1 (b_{r,\gamma})_{\#} \hat{\eta}_n^\pi(s') d\mathcal{P}_R(dr \mid s, a), \quad \forall s \in \mathcal{S},$$

where $(\widehat{Z}_{s,a,s'})_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}}$ are zero-mean Gaussians with

$$\text{Cov}(\widehat{Z}_{s_1, a_1, s'_1}, \widehat{Z}_{s_2, a_2, s'_2}) = \mathbb{1}_{\{(s_1, a_1) = (s_2, a_2)\}} \widehat{P}(s'_1 | s_1, a_1) \left(\mathbb{1}_{\{s'_1 = s'_2\}} - \widehat{P}(s'_2 | s_1, a_1) \right),$$

and

$$\widehat{z}_1(p) := \inf \left\{ t \mid \mathbb{P} \left(\sup_{f \in \mathcal{F}_{W_1}} [(\mathcal{I} - \widehat{\mathcal{T}}^\pi)^{-1} \widehat{\mathbb{G}}^\pi](s) f \leq t \right) \geq p \right\}.$$

Then $\widehat{z}_1(p) \xrightarrow{P} z_1(p)$ if $z_1(\cdot)$ is continuous at p .

Furthermore, if Assumption 1 holds, we have

$$\sup_{f \in \mathcal{F}_{W_1}} [(\mathcal{I} - \widehat{\mathcal{T}}^\pi)^{-1} \widehat{\mathbb{G}}^\pi](s) f = \int_0^{\frac{1}{1-\gamma}} \left| [(\mathcal{I} - \widehat{\mathcal{T}}^\pi)^{-1} \widehat{\mathbb{G}}^\pi](s) \mathbb{1}_{(-\infty, x]} \right| dx,$$

which can be computed efficiently, and $z_1(\cdot)$ is continuous at any $p \in (0, 1)$.

The proof is in Appendix C.

We can also construct confidence sets in the form of KS balls and TV balls for $\eta^\pi(s)$ when Assumption 1 or Assumption 2 holds.

Proposition 4.3. For some fixed policy π and $s \in \mathcal{S}$, define

$$\rho_2(\alpha) := \frac{z_2(1-\alpha)}{\sqrt{n}}, \text{ where } z_2(p) \text{ is defined as the } p\text{-quantile of } \sup_{f \in \mathcal{F}_{\text{KS}}} [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \widetilde{\mathbb{G}}^\pi](s) f,$$

$$\rho_3(\alpha) := \frac{z_3(1-\alpha)}{\sqrt{n}}, \text{ where } z_3(p) \text{ is defined as the } p\text{-quantile of } \sup_{f \in \mathcal{F}_{\text{TV}}} [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \widetilde{\mathbb{G}}^\pi](s) f.$$

Then we have

- (a) $\lim_{n \rightarrow \infty} \mathbb{P}(\text{KS}(\widehat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \rho_2(\alpha)) = 1 - \alpha$ under Assumption 1;
- (b) $\lim_{n \rightarrow \infty} \mathbb{P}(\text{TV}(\widehat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \rho_3(\alpha)) = 1 - \alpha$ under Assumption 2.

$\rho_2(\alpha)$ and $\rho_3(\alpha)$ asymptotically describe the quantile of $\text{KS}(\widehat{\eta}_n^\pi(s), \eta^\pi(s))$, $\text{TV}(\widehat{\eta}_n^\pi(s), \eta^\pi(s))$, respectively. They are determined using results in Theorem 3.4 and the continuous mapping theorem.

Theorem 4.2. Suppose Assumption 1 or Assumption 2 holds. For some fixed policy π and initial state $s \in \mathcal{S}$, define

$$C_2(\alpha) := \left\{ \eta \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right) \mid \text{KS}(\eta, \widehat{\eta}_n^\pi(s)) \leq \rho_2(\alpha) \right\},$$

$$C_3(\alpha) := \left\{ \eta \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right) \mid \text{TV}(\eta, \widehat{\eta}_n^\pi(s)) \leq \rho_3(\alpha) \right\}.$$

Then we have

$$(a) \lim_{n \rightarrow \infty} \mathbb{P}(\eta^\pi(s) \in C_2(\alpha)) = 1 - \alpha.$$

$$(b) \lim_{n \rightarrow \infty} \mathbb{P}(\eta^\pi(s) \in C_3(\alpha)) = 1 - \alpha.$$

Note that $C_2(\alpha)$ and $C_3(\alpha)$ are asymptotically valid confidence sets. Although they rely on the unknown quantile function $z_2(1 - \alpha)$ and $z_3(1 - \alpha)$, they can be consistently estimated using a plug-in approach as the case of $z_1(1 - \alpha)$.

Proposition 4.4. *For any fixed π and $s \in \mathcal{S}$, define*

$$\hat{z}_2(p) := \inf \left\{ t \mid \mathbb{P} \left(\sup_{f \in \mathcal{F}_{\text{KS}}} \left[(\mathcal{I} - \hat{\mathcal{T}}^\pi)^{-1} \hat{\mathbb{G}}^\pi \right] (s) f \leq t \right) \geq p \right\}.$$

$$\hat{z}_3(p) := \inf \left\{ t \mid \mathbb{P} \left(\sup_{f \in \mathcal{F}_{\text{TV}}} \left[(\mathcal{I} - \hat{\mathcal{T}}^\pi)^{-1} \hat{\mathbb{G}}^\pi \right] (s) f \leq t \right) \geq p \right\}.$$

Then $\hat{z}_2(p) \xrightarrow{p} z_2(p)$ if Assumption 1 holds, $\hat{z}_3(p) \xrightarrow{p} z_3(p)$ if Assumption 2 holds.

One may refer to the proof in Appendix C.

4.2 Inference for Hadamard Differentiable Functionals

We consider the problem of statistical inference for $\phi(\eta^\pi(s))$, where $\phi(\cdot): \ell^\infty(\mathcal{F}_{W_1}) \rightarrow \mathbb{R}$ represents a statistical functional. When $\phi(\cdot)$ is Hadamard differentiable, we can determine the limiting distribution of $\sqrt{n}(\phi(\hat{\eta}_n^\pi(s)) - \phi(\eta^\pi(s)))$ using the functional Delta method. Subsequently, we can construct asymptotic confidence sets for $\phi(\eta^\pi(s))$ based on this result.

Theorem 4.3. *For fixed policy π and $s \in \mathcal{S}$, let $\phi(\cdot): \ell^\infty(\mathcal{F}_{W_1}) \rightarrow \mathbb{R}$ be Hadamard differentiable at $\eta^\pi(s)$ tangentially to $\mathbb{D}_0 \subset \ell^\infty(\mathcal{F}_{W_1})$. Suppose $\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \in \mathbb{D}_0$ and define*

$$C_\phi(\alpha) := \left[\phi(\hat{\eta}_n^\pi(s)) + \frac{z_\phi(\alpha/2)}{\sqrt{n}}, \phi(\hat{\eta}_n^\pi(s)) + \frac{z_\phi(1 - \alpha/2)}{\sqrt{n}} \right],$$

where z_ϕ is the quantile function of $\phi'_{\eta^\pi(s)} \left(\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \right)$, which is indeed a one-dimensional gaussian variable with zero means. We have

$$\lim_{n \rightarrow \infty} \mathbb{P}(\phi(\eta^\pi(s)) \in C_\phi(\alpha)) = 1 - \alpha.$$

Under Assumption 1 or Assumption 2, we have similar results for $\phi(\cdot): \ell^\infty(\mathcal{F}_{\text{KS}}) \rightarrow \mathbb{R}$ or $\ell^\infty(\mathcal{F}_{\text{TV}}) \rightarrow \mathbb{R}$ that is Hadamard differentiable at $\eta^\pi(s)$.

Theorem 4.3 follows directly from our asymptotic results described in Theorem 3.4 and functional delta method (Theorem 20.8 in [40]). Since the derivative ϕ' is continuous, the plug-in approach is still valid for estimating z_ϕ .

Proposition 4.5. *Whenever the Hadamard derivative ϕ' is properly defined, let*

$$\hat{z}_\phi(p) := \inf \left\{ t \mid \mathbb{P} \left(\phi'_{\hat{\eta}_n^\pi(s)} \left(\left[(\mathcal{I} - \hat{\mathcal{T}}^\pi)^{-1} \hat{\mathbb{G}}^\pi \right] (s) \right) \leq t \right) \geq p \right\},$$

Then $\hat{z}_\phi(p) \xrightarrow{P} z_\phi(p)$ if $z_\phi(\cdot)$ is continuous at p .

The proof of the proposition above is nearly identical to those of Proposition 4.2 and Proposition 4.4.

We demonstrate the use of Theorem 4.3 with three concrete examples.

Example 4.1 (The r th moment of returns). *We first consider a simple example of statistical inference for r th moments of returns. Let $\phi_r(\mu) := \mathbb{E}_{X \sim \mu}(X^r)$, where μ is a signed measure supported on $\left[0, \frac{1}{1-\gamma}\right]$. It can be easily verified that $\phi(\cdot): \ell^\infty(\mathcal{F}_{W_1}) \rightarrow \mathbb{R}$ is Hadamard differentiable with the derivative $\phi'_r(h) = \phi_r(h) = \mathbb{E}_{X \sim h}(X^r)$ for any signed measure h supported on $\left[0, \frac{1}{1-\gamma}\right]$. Then by Theorem 4.3 we have*

$$\sqrt{n}(\phi_r(\hat{\eta}_n^\pi(s)) - \phi_r(\eta^\pi(s))) \rightsquigarrow \phi_r \left(\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \right),$$

according to which we may perform statistical inference for $\phi(\eta^\pi(s))$. When $r = 1$, we have

$$\sqrt{n}(\hat{V}^\pi(s) - V^\pi(s)) \rightsquigarrow \left[(I - \gamma P^\pi)^{-1} \tilde{G}^\pi \right] (s).$$

Here $P^\pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{S}}$ is the transition matrix under policy π , \hat{V}^π and V^π are the estimated value function and ground-truth value function. \tilde{G}^π is defined as

$$\tilde{G}^\pi(s) := \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} Z_{s,a,s'} V^\pi(s'),$$

where Z is a gaussian vector as defined in Theorem 3.4. This recovers the results of limiting distributions of the errors of model-based policy evaluations when a generative model is available.

Another simple corollary is the limiting distribution of the variance of returns. Concretely, we have

$$\sqrt{n}(\text{Var}_{X \sim \hat{\eta}_n^\pi(s)}(X) - \text{Var}_{X \sim \eta^\pi(s)}(X)) \rightsquigarrow \phi_2 \left(\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \right) - 2\phi_1 \left(\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \right) \phi_1(\eta^\pi(s)).$$

Example 4.2 (Quantiles of returns). We next consider statistical inference for quantiles of returns when Assumption 1 holds. Let $\phi_p(\mu) := \inf \{t \mid \mu \mathbf{1}_{(-\infty, t]} \geq p\}$ be the p -quantile of probability distribution μ . Lemma A.2 in the proof of 2.5 indicates $\eta^\pi(s)$ must have bounded density. Hence we have ϕ_p is Hadamard differentiable tangentially to $C \left[0, \frac{1}{1-\gamma}\right]$ by Lemma 21.4 in [40]. And the derivative $\phi'_p(\eta^\pi(s))$ is the map $h \mapsto -\frac{h(\phi_p(\eta^\pi(s)))}{g(\phi_p(\eta^\pi(s)))}$, where g is the density of $\eta^\pi(s)$. Therefore, the cumulative distribution function of $\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s)$ is in $C \left[0, \frac{1}{1-\gamma}\right]$ almost surely, we have

$$\sqrt{n}(\phi_p(\hat{\eta}_n^\pi(s)) - \phi_p(\eta^\pi(s))) \rightsquigarrow -\frac{\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s) \mathbf{1}_{(-\infty, \phi_p(\eta^\pi(s)))}}{g(\phi_p(\eta^\pi(s)))},$$

which may lead to asymptotically valid inferential procedures for quantiles of returns.

Example 4.3 (Uniform Advantage). Policy improvement [36] is a key ingredient of many reinforcement learning algorithms. The goal is to find a new policy π such that the advantage function $V^\pi(s_0) - V^{\pi_0}(s_0) \geq 0$ where π_0 is a given baseline policy. Here we propose a new notion of policy improvement called (near)-uniform policy improvement. Specifically, the aim is to find a new policy π such that the uniform advantage $\phi(\eta^\pi(s_0), \eta^{\pi_0}(s_0)) := \mathbb{P}(G^\pi(s_0) \geq G^{\pi_0}(s_0))$ is above some threshold.

A natural estimator of $\phi(\eta^\pi(s_0), \eta^{\pi_0}(s_0))$ is $\phi(\hat{\eta}_n^\pi(s_0), \hat{\eta}_n^{\pi_0}(s_0))$. For technical convenience we assume $\hat{\eta}_n^\pi(s_0)$ and $\hat{\eta}_n^{\pi_0}(s_0)$ are estimated using two different dataset. From Lemma 20.10 in [40], ϕ is Hadamard differentiable tangentially to $C[0, 1/(1-\gamma)]$. For $h_1, h_2 \in C[0, 1/(1-\gamma)]$, the derivative $\phi'(\eta^\pi(s_0), \eta^{\pi_0}(s_0))$ is $(h_1, h_2) \mapsto h_2 - \eta^\pi(s_0)h_2 + \eta^{\pi_0}(s_0)h_1$. When Assumption 1 is true,

$$\sqrt{n}[(\hat{\eta}_n^\pi(s_0), \hat{\eta}_n^{\pi_0}(s_0)) - (\eta^\pi(s_0), \eta^{\pi_0}(s_0))] \rightsquigarrow \left(\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s_0), \left[(\mathcal{I} - \mathcal{T}^{\pi_0})^{-1} \tilde{\mathbb{G}}^{\pi_0} \right] (s_0) \right)$$

and the cumulative distribution functions of $\left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right] (s_0)$ and $\left[(\mathcal{I} - \mathcal{T}^{\pi_0})^{-1} \tilde{\mathbb{G}}^{\pi_0} \right] (s_0)$ (denoted as F_1 and F_2) are in $C[0, 1/(1-\gamma)]$ almost surely. Thus we have

$$\sqrt{n}[\phi(\hat{\eta}_n^\pi(s_0), \hat{\eta}_n^{\pi_0}(s_0)) - \phi(\eta^\pi(s_0), \eta^{\pi_0}(s_0))] \rightsquigarrow F_2 - \eta^\pi(s_0)F_2 + \eta^{\pi_0}(s_0)F_1.$$

5 Numerical Simulations

In this section, we conduct numerical simulations to validate our theoretical findings as well as the proposed inferential procedures. All of the numerical simulations are conducted on a desktop computer with a single TITAN RTX GPU. The code is available in <https://github.com/zhangliangyu32/EstimationAndInferenceDistributionalRL>.

5.1 Implementations

To make computations tractable, we confine the return distributions to the class of categorical distributions. A vector of categorical distributions is defined as $\eta = (\eta(s_1), \dots, \eta(s_{|S|}))$, where $\eta(s) := \sum_{k=0}^K w_k \delta_{x_k}$, with weights $\sum_{k=0}^K w_k = 1$ and particles $x_k := \frac{k}{(K+1)(1-\gamma)}$. We set $K = 1000$, which is large enough to make the categorical class rich enough and good approximations of continuous return distributions.

The categorical distributions can be updated with a categorical version of distributional dynamical programming [1], which is also a good approximation of the original version of distributional dynamic programming considered in our paper when K is large. Throughout our simulation studies, the ground-truth return distributions η^π are obtained via a sufficiently large number of iterations of distributional dynamic programming with the ground-truth distributional Bellman operator \mathcal{T}^π . The estimated return distributions $\hat{\eta}_n^\pi$ are obtained by the same procedure except for the ground-truth distributional Bellman operator \mathcal{T}^π is replaced by the estimated distributional Bellman operator $\hat{\mathcal{T}}_n^\pi$.

Another issue that may cause computational intractability is that in our inferential procedures, we must explicitly form the operator $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$. We instead use a truncated Neumann series to approximate $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$, *i.e.* $(\mathcal{I} - \mathcal{T}^\pi)^{-1} \approx \sum_{j=0}^J (\mathcal{T}^\pi)^j$ with J sufficiently large. In summary, by these approximation techniques, our implementations achieve computational tractability while ensuring that the approximation error is negligible compared to the statistical error that is of primary interest.

5.2 Finite-sample Convergence Performance

We first investigate the finite-sample convergence performances of empirical distributional dynamic programming and verify our non-asymptotic results. We perform the simulations in randomly-generated tabular MDPs with $|\mathcal{S}| = 5$, $|\mathcal{A}| = 2$ and $\gamma \in \{0.7, 0.8, 0.9, 0.97\}$. WLOG, we always use

the first state s_1 as the initial state. The reward distribution is chosen to be truncated gaussians. Specifically, $\mathcal{P}_R(\cdot | s, a)$ is set to be $\mathcal{N}(l_{s,a}, 0.1)$ truncated to $[0, 1]$, with the location parameter $l_{s,a}$ randomly determined. The dataset we use to form the estimator $\hat{\eta}_n^\pi$ is obtained via a generative model with $n \in \{10, 100, 1000, 10000\}$. We repeat the estimation process for 100 times and report the averaged results. The numerical results are displayed in figures listed as follows.

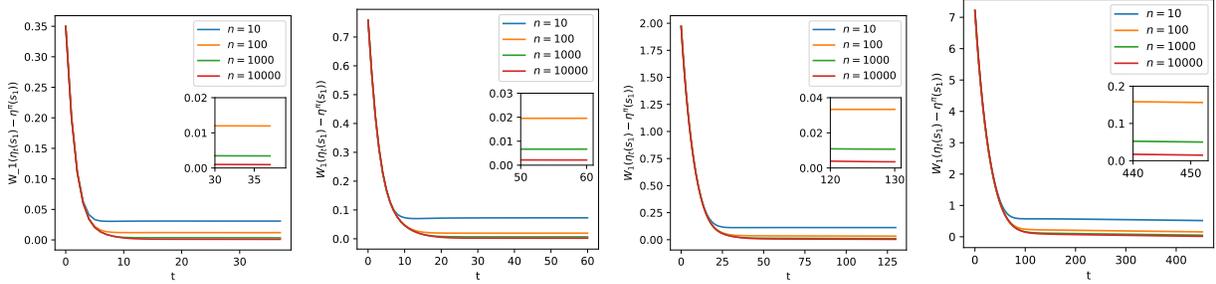


Figure 2. Two-phase convergence of $W_1(\eta^t(s_1), \eta^\pi(s_1))$ with different sample size. t is the iteration number. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

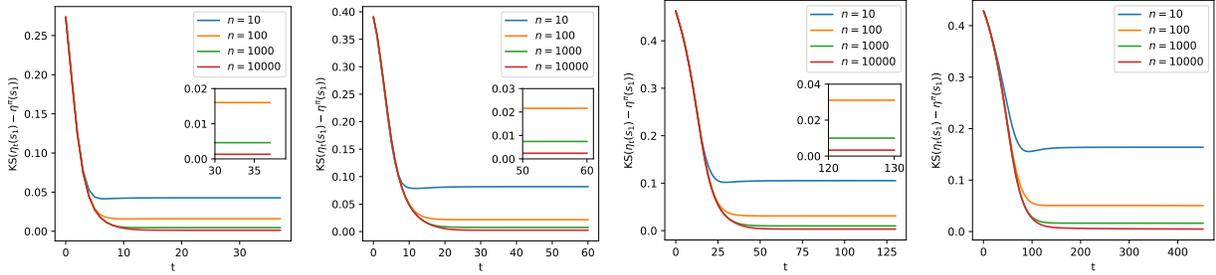


Figure 3. Two-phase convergence of $KS(\eta^t(s_1), \eta^\pi(s_1))$ with different sample size. t is the iteration number. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

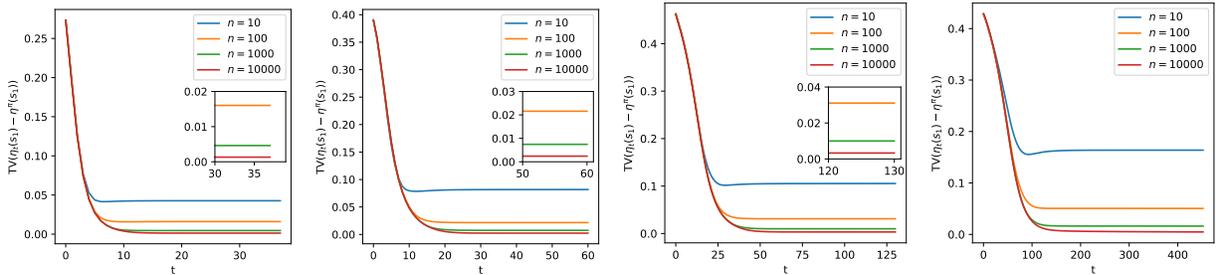


Figure 4. Two-phase convergence of $TV(\eta^t(s_1), \eta^\pi(s_1))$ with different sample size. t is the iteration number. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

In Figure 2-4, we show the convergence performance of empirical distributional dynamic programming measured by W_1 metric, KS distance and TV distance, respectively. Note that in all cases the convergence consists of two phases. In the first phase, the dynamic programming algorithm does not converge and we may observe a linear convergence rate depicted in Proposition 2.4, 2.5, and 2.7.

In the second phase, the error terms are dominated by the statistical error, *i.e.* $W_1(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$, $\text{KS}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ or $\text{TV}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$, which exhibits strong correlations with n and $1/(1-\gamma)$.

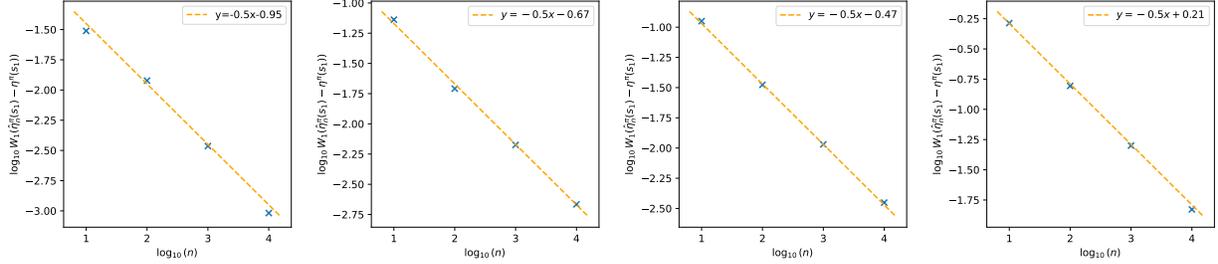


Figure 5. Two-phase convergence of $W_1(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ with different sample size. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

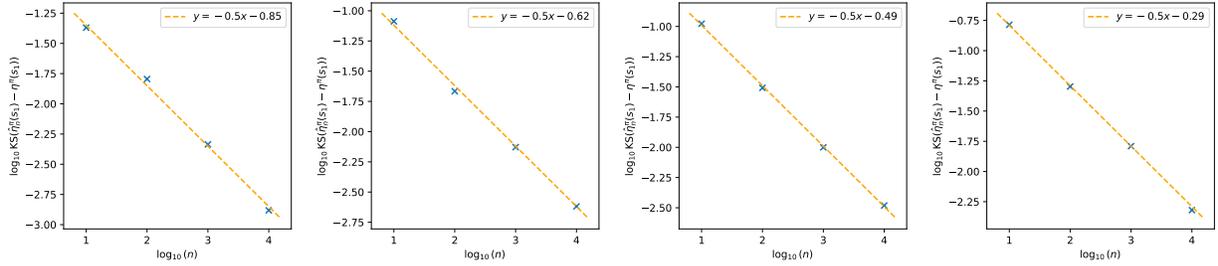


Figure 6. Two-phase convergence of $\text{KS}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ with different sample size. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

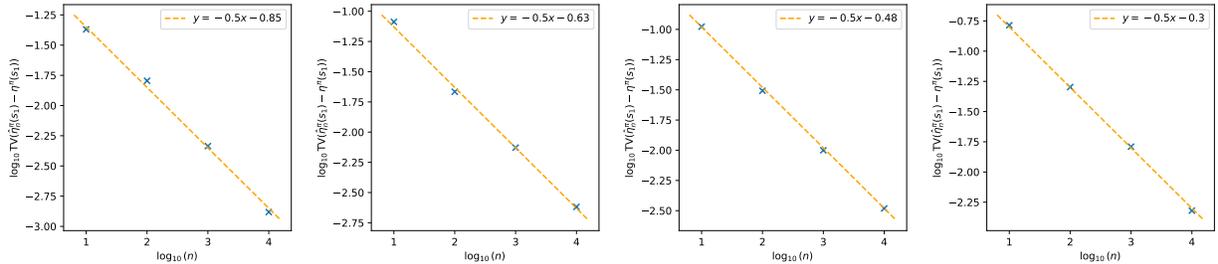


Figure 7. Two-phase convergence of $\text{TV}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ with different sample size. From left to right: $\gamma = 0.7$; $\gamma = 0.8$; $\gamma = 0.9$; $\gamma = 0.97$.

We also try to examine how the error terms $W_1(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$, $\text{KS}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ or $\text{TV}(\hat{\eta}_n^\pi(s_1), \eta^\pi(s_1))$ change as n scales up. In Figure 5-7, we may verify that for all cases, the convergence rates are indeed of the typical $n^{-\frac{1}{2}}$ order as described in Theorem 3.1, 3.2 and 3.3.

5.3 Validity of Inferential Procedures

We also perform numerical simulations to validate the inferential procedures proposed in Section 4. The environment is exactly the same as that of the previous section with γ fixed as 0.9. The confidence sets are constructed using plug-in approaches. The nominal coverage probability is

Type of Confidence Sets	W_1 ball		KS ball		TV ball	
n	CR	CSR \pm std	CR	CSR \pm std	CR	CSR \pm std
5	0.918	0.3467 \pm 0.0719	0.939	0.3302 \pm 0.0538	0.934	0.3310 \pm 0.0530
10	0.945	0.2528 \pm 0.0366	0.934	0.2364 \pm 0.0250	0.956	0.2367 \pm 0.0246
100	0.941	0.0804 \pm 0.0041	0.956	0.0759 \pm 0.0032	0.944	0.0761 \pm 0.0030
1000	0.945	0.0255 \pm 0.0008	0.951	0.0241 \pm 0.0007	0.950	0.0241 \pm 0.0007

Table 1. Coverage rate (CR) and confidence set radius (CSR) of our proposed non-parametric confidence sets for $\eta^\pi(s_1)$ under different choices of n .

Functionals of Interest	variance		0.1 quantile		0.9 quantile	
n	CR	CSR \pm std	CR	CSR \pm std	CR	CSR \pm std
5	0.928	0.0627 \pm 0.0202	0.917	0.4020 \pm 0.0929	0.916	0.3210 \pm 0.0660
10	0.939	0.0442 \pm 0.0101	0.945	0.2949 \pm 0.0436	0.930	0.2288 \pm 0.0308
100	0.959	0.0137 \pm 0.0010	0.949	0.0912 \pm 0.0050	0.946	0.0733 \pm 0.0036
1000	0.946	0.0043 \pm 0.0002	0.935	0.0289 \pm 0.0010	0.952	0.0232 \pm 0.0008

Table 2. Coverage rate (CR) and confidence set radius (CSR) of our proposed confidence intervals for different Hardamard differentiable statistical functionals of $\eta^\pi(s_1)$ under different choices of n .

set to be 0.95, the quantiles of the estimated limiting distributions are computed using Monte Carlo methods. Here our implementations of Monte Carlos are fully vectorized to further improve computational efficiency. We repeat our inferential procedures for 1000 times and report the empirical coverage rate and averaged radius of confidence sets. The results are presented in the following tables. We observe the empirical coverage rates approach the nominal confidence level and the radius of confidence sets decreases as n increases in all cases.

6 Proof Outlines

Before presenting our proofs, we first define some notations. For any signed Borel measure μ , we define $\|\mu\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} |\mu f|$, with \mathcal{F} being some function class \mathcal{F} supported on $\left[0, \frac{1}{1-\gamma}\right]$. We use $M_{\mathcal{F}}$ to denote the vector space of signed Borel measures with finite $\|\cdot\|_{\mathcal{F}}$ norm and zero measure, formally, define \mathcal{B}_0 as the Borel sets in $\left[0, \frac{1}{1-\gamma}\right]$,

$$M_{\mathcal{F}}^0 := \left\{ \mu \text{ signed measure on } \left(\left[0, \frac{1}{1-\gamma}\right], \mathcal{B}_0 \right) \mid \mu \left(\left[0, \frac{1}{1-\gamma}\right] \right) = 0, \|\mu\|_{\mathcal{F}} < \infty \right\}.$$

Let $\ell^\infty(\mathcal{F})$ be the space of bounded real-valued functions on \mathcal{F} . $\ell^\infty(\mathcal{F})$ is a Banach space if we equip it with the supreme norm $\|\cdot\|_{\mathcal{F}}$. Note that $(M_{\mathcal{F}}^0, \|\cdot\|_{\mathcal{F}})$ can be embedded into $(\ell^\infty(\mathcal{F}), \|\cdot\|_{\mathcal{F}})$ by mapping μ to the process $f \mapsto \mu f$.

Recall that the W_1 metric, KS distance, and the TV are all integral probability metrics. Therefore,

if $\mu \in M_{\mathcal{F}}^0$ can be written as the difference of two probability distributions $\mu_+, \mu_- \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)$, we have

- $\|\mu\|_{\mathcal{F}_{W_1}} = W_1(\mu_+, \mu_-)$, where $\mathcal{F}_{W_1} := \{f \mid f \text{ is 1-Lipschitz}\}$.
- $\|\mu\|_{\mathcal{F}_{\text{KS}}} = \text{KS}(\mu_+, \mu_-)$, where $\mathcal{F}_{\text{KS}} := \{\mathbf{1}\{\cdot \leq z\} \mid z \in \mathbb{R}\}$.
- $\|\mu\|_{\mathcal{F}_{\text{TV}}} = \text{TV}(\mu_+, \mu_-)$, where $\mathcal{F}_{\text{TV}} := \{f \mid f \text{ is measurable and } \|f\|_{\infty} \leq 1\}$.

The distributional Bellman operator \mathcal{T}^π can be naturally extended to $(M_{\mathcal{F}}^0)^{\mathcal{S}}$ without modifying its original definition. Here, the product $(M_{\mathcal{F}}^0)^{\mathcal{S}}$ also constitutes a normed vector space, where the norm is defined as $\sup_{s \in \mathcal{S}} \|\cdot\|_{\mathcal{F}}$.

Proposition 6.1. \mathcal{T}^π is a linear operator on $(M_{\mathcal{F}}^0)^{\mathcal{S}}$.

The proof is straightforward by noting the linearity of the push-forward operation.

6.1 Analysis of Theorem 3.1

Let

$$M_{\mathcal{F}_{W_1}}^0 := \left\{ \mu \text{ signed measure on } \left(\left[0, \frac{1}{1-\gamma}\right], \mathcal{B}_0 \right) \mid \|\mu\|_{\mathcal{F}_{W_1}} < \infty, \mu \left(\left[0, \frac{1}{1-\gamma}\right] \right) = 0 \right\},$$

for any vector of probability measures $\mu, \nu \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)^{\mathcal{S}}$, we always have $\mu - \nu \in (M_{\mathcal{F}_{W_1}}^0)^{\mathcal{S}}$. Therefore, we have

$$\begin{aligned} \hat{\eta}_n^\pi - \eta^\pi &= \widehat{\mathcal{T}}_n^\pi \hat{\eta}_n^\pi - \mathcal{T}^\pi \eta^\pi \\ &= \widehat{\mathcal{T}}_n^\pi \hat{\eta}_n^\pi - \widehat{\mathcal{T}}_n^\pi \eta^\pi + \widehat{\mathcal{T}}_n^\pi \eta^\pi - \mathcal{T}^\pi \eta^\pi \\ &= \widehat{\mathcal{T}}_n^\pi (\hat{\eta}_n^\pi - \eta^\pi) + (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi. \end{aligned}$$

Rearranging terms yields

$$\left(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi) = \left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi.$$

We first investigate the concentration behavior of $\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi$. Formally, we have

Lemma 6.1. For any fixed policy π , we have

$$\mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} \leq \sqrt{\frac{9 \log |\mathcal{S}|}{n(1-\gamma)^2}}.$$

And for any $\delta \in (0, 1)$,

$$\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} \leq \frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2}}{\sqrt{n(1-\gamma)^2}}$$

with probability greater than $1 - \delta$.

The high-level idea of proof is to first study the concentration of $\left| \widehat{F}_s(x) - F_s(x) \right|$, where $\widehat{F}_s(x)$ is defined to be the cumulative distribution function of $\left[\widehat{\mathcal{T}}_n^\pi \eta^\pi \right] (s)$ and $F_s(x)$ is defined as the cumulative distribution function of $[\mathcal{T}^\pi \eta^\pi] (s)$. And then use the fact that

$$\left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} = \int_0^{\frac{1}{1-\gamma}} \left| \widehat{F}_s(x) - F_s(x) \right| dx.$$

The full proof is in Appendix D.

The next step is to relate the error term $\widehat{\eta}_n^\pi - \eta^\pi$ with $\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi$. Since

$$\left(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi \right) \left(\widehat{\eta}_n^\pi - \eta^\pi \right) = \left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi,$$

this can be immediately accomplished if $\left(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi \right)$ is invertible on $\left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$. The invertibility does not hold in general. However, we find $\left(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi \right)$ is invertible on the closure $\left(\overline{M_{\mathcal{F}_{W_1}}^0} \right)^\mathcal{S}$, which suffices for our analysis.

Lemma 6.2. *The operator $\left(\mathcal{I} - \mathcal{T}^\pi \right)$ is invertible on $\left(\overline{M_{\mathcal{F}_{W_1}}^0} \right)^\mathcal{S}$ and $\left(\mathcal{I} - \mathcal{T}^\pi \right)^{-1} = \sum_{i=0}^{\infty} \left(\mathcal{T}^\pi \right)^i$. Also, the operator norm of $\left(\mathcal{I} - \mathcal{T}^\pi \right)^{-1}$ is upper bounded by $\frac{1}{1-\gamma}$. Here \mathcal{T}^π can be replaced by any valid distributional Bellman operator.*

Lemma 6.2 not only shows $\left(\mathcal{I} - \mathcal{T}^\pi \right)$ is invertible but constructs the inverse explicitly, thereby facilitating computational convenience. Note that for $\mu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$, it is not necessarily true that $\left(\mathcal{I} - \mathcal{T}^\pi \right)^{-1} \mu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$ because the space $\left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$ is not complete. Precisely, we have $\left(\mathcal{I} - \mathcal{T}^\pi \right)^{-1} \mu \in \left(\overline{M_{\mathcal{F}_{W_1}}^0} \right)^\mathcal{S}$ that is a closed subspace of $\left(\ell^\infty(\mathcal{F}_{W_1}) \right)^\mathcal{S}$ for any $\mu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$.

To prove Lemma 6.2, the main idea is to show the convergence of the Neumann series $\sum_{i=1}^{\infty} \left(\mathcal{T}^\pi \right)^i$. For any $\mu \in M_{\mathcal{F}_{W_1}}^0$, one has the Jordan decomposition $\mu = a_\mu(\mu_+ - \mu_-)$ such that a_μ is a positive constant and μ_+, μ_- are probability measures. Thus $\|\mu\|_{\mathcal{F}_{W_1}} = a_\mu W_1(\mu_+, \mu_-)$ and the convergence of Neumann series can be shown utilizing the contraction property of the distributional Bellman operator \mathcal{T}^π . One may refer to Appendix D for detailed proof.

Applying Lemma 6.2 to $\widehat{\mathcal{T}}_n^\pi$ leads to

$$\begin{aligned} \sup_{s \in \mathcal{S}} W_1(\widehat{\eta}_n^\pi(s), \eta^\pi(s)) &= \sup_{s \in \mathcal{S}} \|\widehat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{W_1}} \\ &= \sup_{s \in \mathcal{S}} \left\| \left[(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)^{-1} (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} \\ &\leq \frac{1}{1 - \gamma} \sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}}, \end{aligned}$$

which finishes our analysis of the W_1 -metric case.

6.2 Analysis of Theorem 3.2

We first give a stronger notion of concentration of $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$.

Lemma 6.3. *Suppose Assumption 1 is true. For any fixed policy π , we have*

$$\mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \leq \frac{3C}{2} \sqrt{\frac{\log |\mathcal{S}|}{n(1 - \gamma)^2}}.$$

And for any $\delta \in (0, 1)$,

$$\sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \leq \frac{C \left(\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2} \right)}{2\sqrt{n(1 - \gamma)^2}}$$

with probability greater than $1 - \delta$.

Different from the proof of Lemma 6.1, this time we first bound the term $|\widehat{p}_s(x) - p_s(x)|$, where $\widehat{p}_s(x)$, $p_s(x)$ are defined to be the density functions of $\left[\widehat{\mathcal{T}}_n^\pi \eta^\pi \right] (s)$, $[\mathcal{T}^\pi \eta^\pi]$ respectively. Then we may draw the conclusion noting that

$$\left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} = \frac{1}{2} \int_0^{\frac{1}{1-\gamma}} |\widehat{p}_s(x) - p_s(x)| dx.$$

The full proof is in Appendix D.

As before, the next step is to show $(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)$ is invertible on some space containing the signed measures of interest. Specifically, let

$$M_{\mathcal{F}_{\text{KS}}}^0 := \left\{ \mu \text{ signed measure on } \left(\left[0, \frac{1}{1-\gamma} \right], \mathcal{B}_0 \right) \mid \|\mu\|_{\mathcal{F}_{\text{KS}}} < \infty, \|\mu\|_{\text{loc}} < \infty, \mu \left(\left[0, \frac{1}{1-\gamma} \right] \right) = 0. \right\}$$

Here $\|\mu\|_{\text{loc}} < \infty$ represents that μ has a density f such that for any $B \in \mathcal{B}_0$, $\mu(B) = \int_B f(x) dx$

and $\|\mu\|_{\text{loc}} := \sup_{x \in [0, 1/(1-\gamma)]} |f(x)|$. When $\mu \in M_{\mathcal{F}_{\text{KS}}}^0$, we can control $\|\mu\|_{\mathcal{F}_{\text{KS}}}$ with $\|\mu\|_{\mathcal{F}_{W_1}}$ and $\|\mu\|_{\text{loc}}$. Formally, we have the following proposition as a generalization of Proposition 2.2.

Proposition 6.2. *Suppose $\mu \in M_{\mathcal{F}_{\text{KS}}}^0$, then $\|\mu\|_{\mathcal{F}_{\text{KS}}} \leq \sqrt{2 \|\mu\|_{\text{loc}} \|\mu\|_{\mathcal{F}_{W_1}}}$.*

The proof is in Appendix D. When Assumption 1 is true, we always have $(\eta^\pi - \hat{\eta}_n^\pi) \in \left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$ as $\sup_{s \in \mathcal{S}} \|\eta^\pi(s) - \hat{\eta}_n^\pi(s)\|_{\text{loc}} \leq C$. Also, if $\mu \in \left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$, then $\mathcal{T}^\pi \mu \in \left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$ (Lemma A.1). Here \mathcal{T}^π can be replaced by any valid distributional Bellman operator, for example, $\hat{\mathcal{T}}_n^\pi$.

Lemma 6.4. *For any valid distributional Bellman operator \mathcal{T}^π , the operator $(\mathcal{I} - \mathcal{T}^\pi)$ is invertible on $\left(\overline{M_{\mathcal{F}_{\text{KS}}}^0}\right)^{\mathcal{S}}$ and $(\mathcal{I} - \mathcal{T}^\pi)^{-1} = \sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$.*

The formal proof can be found in Appendix D. Like the case of Lemma 6.2, for $\mu \in \left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$, $(\mathcal{I} - \mathcal{T}^\pi)^{-1} \mu$ does not necessarily lie in $\left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$. Instead, we have $(\mathcal{I} - \mathcal{T}^\pi)^{-1} \mu \in \left(\overline{M_{\mathcal{F}_{\text{KS}}}^0}\right)^{\mathcal{S}}$ that is a closed subspace of $(\ell^\infty(\mathcal{F}_{\text{KS}}))^{\mathcal{S}}$ for any $\mu \in \left(M_{\mathcal{F}_{\text{KS}}}^0\right)^{\mathcal{S}}$.

Since the inverse $(\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1}$ can be unbounded in $M_{\mathcal{F}_{\text{KS}}}^0$, the analysis is more involved. Here we detour the technical problem with an ‘‘expansion trick’’. For any $s \in \mathcal{S}$, we have

$$\begin{aligned}
& \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} \\
&= \left\| \left[(\mathcal{I} - \hat{\mathcal{T}}_n^\pi)^{-1} (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&= \left\| \left[\sum_{i=0}^{\infty} (\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&\leq \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} + \sum_{i=1}^{\infty} \left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&\leq \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} + \sum_{i=1}^{\infty} \sqrt{2 \left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} \left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\text{loc}}} \\
&\leq \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} + \sum_{i=1}^{\infty} \sqrt{2 \gamma^i \sup_{s' \in \mathcal{S}} \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}} \left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\text{loc}}} \\
&\leq \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \sqrt{2 \sup_{s' \in \mathcal{S}} \left\| \left[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}} \sum_{i=1}^{\infty} \sqrt{\gamma^i} \left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\text{loc}}}.
\end{aligned}$$

Here the second inequality is due to Proposition 6.2, the third inequality is by the contraction property of $\hat{\mathcal{T}}_n^\pi$, and the last inequality holds from the fact $\|\cdot\|_{\mathcal{F}_{\text{KS}}} \leq \|\cdot\|_{\mathcal{F}_{\text{TV}}}$.

Now we need an upper bound for $\left\| \left[(\hat{\mathcal{T}}_n^\pi)^i (\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\text{loc}}$. Assumption 1 implies

$\left\| \left[\left(\widehat{\mathcal{T}}_n^\pi \right)^i \left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\text{loc}} \leq C$, but we can do better here.

Lemma 6.5. *Suppose Assumption 1 holds true. For any $i \geq 1$, $s \in \mathcal{S}$,*

$$\left\| \left[\left(\widehat{\mathcal{T}}_n^\pi \right)^i \left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\text{loc}} \leq C \sup_{s' \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{\text{TV}}}.$$

We defer the proof to Appendix D. The main idea is we first normalize $\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi$ to have a proper Jordan decomposition and then use the fact that under Assumption 1, for any $\nu \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right)^{\mathcal{S}}$, $[\mathcal{T}^\pi \nu](s)$ must have a density function bounded by C as long as $\nu(s)$ has a density, $\forall s \in \mathcal{S}$. Hence the condition $i \geq 1$ is necessary and that is the reason why we break the summation into two parts: $i = 0$ and $i \geq 1$.

To sum up,

$$\begin{aligned} & \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} \\ & \leq \frac{\sqrt{\gamma}}{1-\sqrt{\gamma}} \sqrt{2C \sup_{s' \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{\text{TV}}} \sup_{s' \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}} + \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}. \end{aligned}$$

Combining Lemma 6.1 and Lemma 6.3, we have for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} \leq \frac{C}{2} \left(\frac{\sqrt{2\gamma}}{1-\sqrt{\gamma}} + 1 \right) \frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(2/\delta)}/2}{\sqrt{n(1-\gamma)^2}}$$

Note that $\frac{1}{1-\sqrt{\gamma}} = \frac{1+\gamma}{1-\gamma} \leq \frac{2}{1-\gamma}$ when $\gamma \in (0, 1)$, thus

$$\sup_{s \in \mathcal{S}} \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} \leq \frac{C' \left(\sqrt{\log |\mathcal{S}|} + \sqrt{\log(1/\delta)} \right)}{\sqrt{n(1-\gamma)^4}},$$

where C' is some constant depending on C in Assumption 1. We also have

$$\begin{aligned} \mathbb{E} \sup_{s \in \mathcal{S}} \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} & \leq C' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}} + \int_0^\infty \mathbb{P} \left(\sup_{s \in \mathcal{S}} \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{KS}}} > C' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}} + t \right) dt \\ & \leq C' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}} + \int_0^\infty \exp \left(-\frac{n(1-\gamma)^4 t^2}{C'^2} \right) dt \\ & \leq C'' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}}, \end{aligned}$$

where C'' is some constant depending on C in Assumption 1.

6.3 Analysis of Theorem 3.3

According to Proposition 2.6, we also have the following concentration results of $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$ as a corollary of Lemma 6.3.

Corollary 6.1. *Suppose Assumption 2 is true. For any fixed policy π , we have*

$$\mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \leq 3M \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^2}}.$$

And for any $\delta \in (0, 1)$,

$$\sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \leq \frac{M \left(\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2} \right)}{\sqrt{n(1-\gamma)^2}}$$

with probability greater than $1 - \delta$.

We consider the following space containing the signed measures of interest. Specifically, let

$$M_{\mathcal{F}_{\text{TV}}}^0 := \left\{ \mu \text{ signed measure on } \left(\left[0, \frac{1}{1-\gamma} \right], \mathcal{B}_0 \right) \mid \|\mu\|_{\mathcal{F}_{\text{TV}}} < \infty, \|\mu\|_{H_1^1} < \infty, \mu \left(\left[0, \frac{1}{1-\gamma} \right] \right) = 0. \right\}$$

Here we abuse the notation to define $\|\mu\|_{H_1^1} = \|p\|_{H_1^1}$, where f is the density of μ . When $\mu \in M_{\mathcal{F}_{\text{TV}}}^0$, we can control $\|\mu\|_{\mathcal{F}_{\text{TV}}}$ with $\|\mu\|_{\mathcal{F}_{W_1}}$ and $\|\mu\|_{H_1^1}$. Formally, we have the following proposition as a generalization of Proposition 2.3.

Proposition 6.3. *Suppose $\mu \in M_{\mathcal{F}_{\text{TV}}}^0$, then $\|\mu\|_{\mathcal{F}_{\text{TV}}} \leq \sqrt{K \|\mu\|_{H_1^1} \|\mu\|_{\mathcal{F}_{W_1}}}$.*

The proof is in Appendix D. When Assumption 2 is true, we always have $(\eta^\pi - \widehat{\eta}_n^\pi) \in \left(M_{\mathcal{F}_{\text{TV}}}^0 \right)^\mathcal{S}$ (Lemma A.4). Also, if $\mu \in \left(M_{\mathcal{F}_{\text{TV}}}^0 \right)^\mathcal{S}$, then $\mathcal{T}^\pi \mu \in \left(M_{\mathcal{F}_{\text{TV}}}^0 \right)^\mathcal{S}$ (A.3). Here \mathcal{T}^π can be replaced by any valid distributional Bellman operator, for example, $\widehat{\mathcal{T}}_n^\pi$.

Lemma 6.6. *For any valid distributional Bellman operator \mathcal{T}^π , the operator $(\mathcal{I} - \mathcal{T}^\pi)$ is invertible on $\left(\overline{M_{\mathcal{F}_{\text{TV}}}^0} \right)^\mathcal{S}$ and $(\mathcal{I} - \mathcal{T}^\pi)^{-1} = \sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$.*

The full proof is in Appendix D. We have $(\mathcal{I} - \mathcal{T}^\pi)^{-1} \mu \in \left(\overline{M_{\mathcal{F}_{\text{TV}}}^0} \right)^\mathcal{S}$ that is a closed subspace of $(\ell^\infty(\mathcal{F}_{\text{TV}}))^\mathcal{S}$ for any $\mu \in \left(M_{\mathcal{F}_{\text{TV}}}^0 \right)^\mathcal{S}$.

Similar to the analysis of Theorem 3.2, we also use an ‘‘expansion trick’’ to tackle the unbound-

edness issue of $(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)^{-1}$. For any $s \in \mathcal{S}$,

$$\begin{aligned}
& \|\widehat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{TV}}} \\
&= \left\| \left[(\mathcal{I} - \widehat{\mathcal{T}}_n^\pi)^{-1} (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \\
&= \left\| \left[\sum_{i=0}^{\infty} (\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \\
&\leq \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \sum_{i=1}^{\infty} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \\
&\leq \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \sum_{i=1}^{\infty} \sqrt{K} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{H_1^1} \\
&\leq \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \sum_{i=1}^{\infty} \sqrt{K \gamma^i \sup_{s' \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{H_1^1} \\
&\leq \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \sqrt{K \sup_{s' \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}}} \sum_{i=1}^{\infty} \sqrt{\gamma^i} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{H_1^1}.
\end{aligned}$$

The second inequality is due to Proposition 6.3, and the third inequality is by the contraction property of $\widehat{\mathcal{T}}_n^\pi$.

Next, we bound the term $\left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{H_1^1}$.

Lemma 6.7. *Suppose Assumption 2 holds true. For any $i \geq 1$, $s \in \mathcal{S}$,*

$$\left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{H_1^1} \leq 2M \sup_{s \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}.$$

One can refer to Appendix D for proof. Like the proof of Lemma 6.5, we also deploy the normalization technique. And the condition $i \geq 1$ is also necessary.

Putting the pieces together, we have

$$\begin{aligned}
& \|\widehat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{TV}}} \\
&\leq \frac{\sqrt{2KM\gamma}}{1 - \sqrt{\gamma}} \sqrt{\sup_{s' \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{\text{TV}}} \sup_{s' \in \mathcal{S}} \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s') \right\|_{\mathcal{F}_{W_1}}} + \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}.
\end{aligned}$$

According to Lemma 6.1 and Corollary 6.1,

$$\begin{aligned}\|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{TV}}} &\leq M \left(\frac{\sqrt{2K\gamma}}{1 - \sqrt{\gamma}} + 1 \right) \frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(2/\delta)/2}}{\sqrt{n(1-\gamma)^2}} \\ &\leq \frac{K' \left(\sqrt{\log |\mathcal{S}|} + \sqrt{\log(1/\delta)} \right)}{\sqrt{n(1-\gamma)^4}}\end{aligned}$$

with probability at least $1 - \delta$, $\forall \delta \in (0, 1)$. K' is an absolute constant depending only on M in Assumption 2. Besides,

$$\begin{aligned}\mathbb{E} \sup_{s \in \mathcal{S}} \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{TV}}} &\leq K' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}} + \int_0^\infty \mathbb{P} \left(\sup_{s \in \mathcal{S}} \|\hat{\eta}_n^\pi(s) - \eta^\pi(s)\|_{\mathcal{F}_{\text{TV}}} > K' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}} + t \right) dt \\ &\leq K'' \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^4}},\end{aligned}$$

where K'' is some constant depending on M in Assumption 2.

6.3.1 Analysis of Theorem 3.4

Weak Convergence in $\ell^\infty(\mathcal{F}_{W_1})$ We have

$$\begin{aligned}\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi) &= \sqrt{n} \left(\hat{\mathcal{T}}_n^\pi \hat{\eta}_n^\pi - \mathcal{T}^\pi \eta^\pi \right) \\ &= \sqrt{n} \left(\hat{\mathcal{T}}_n^\pi \hat{\eta}_n^\pi - \hat{\mathcal{T}}_n^\pi \eta^\pi + \hat{\mathcal{T}}_n^\pi \eta^\pi - \mathcal{T}^\pi \eta^\pi \right) \\ &= \sqrt{n} \hat{\mathcal{T}}_n^\pi (\hat{\eta}_n^\pi - \eta^\pi) + \sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \\ &= \sqrt{n} \mathcal{T}^\pi (\hat{\eta}_n^\pi - \eta^\pi) + \sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi + \sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi).\end{aligned}$$

Rearranging terms yields

$$\sqrt{n}(\mathcal{I} - \mathcal{T}^\pi) (\hat{\eta}_n^\pi - \eta^\pi) = \underbrace{\sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi}_{(1)} + \underbrace{\sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi)}_{(2)}.$$

Both term (1) and term (2) are in $\left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$. Next, we can show that term (1) converges weakly to a mixture of probability distributions and term (2) is negligible. One may refer to Appendix D for detailed proof.

Lemma 6.8. *For any $s \in \mathcal{S}$, $\sqrt{n} \left[\left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s)$ converge weakly to the process $f \mapsto \tilde{\mathbb{G}}^\pi(s) f$*

in $\ell^\infty(\mathcal{F}_{W_1})$. Here the random element $\tilde{\mathbb{G}}^\pi$ is defined as

$$\tilde{\mathbb{G}}^\pi(s) := \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} Z_{s,a,s'} \int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \quad \forall s \in \mathcal{S},$$

where $(Z_{s,a,s'})_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}}$ are zero-mean Gaussians with

$$\text{Cov}(Z_{s_1,a_1,s'_1}, Z_{s_2,a_2,s'_2}) = \mathbb{1}\{(s_1, a_1) = (s_2, a_2)\} P(s'_1 | s_1, a_1) (\mathbb{1}\{s'_1 = s'_2\} - P(s'_2 | s_1, a_1)).$$

Lemma 6.9. For any $s \in \mathcal{S}$, we have $\left\| \sqrt{n} \left[\left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi) \right] (s) \right\|_{\mathcal{F}_{W_1}} = o_P(1)$.

Recall that we have previously demonstrated $\sqrt{n}(\mathcal{I} - \mathcal{T}^\pi)(\hat{\eta}_n^\pi - \eta^\pi) \rightsquigarrow \tilde{\mathbb{G}}^\pi$. Thus our final step is to establish $\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi) \rightsquigarrow (\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi$. This step can be accomplished by continuous mapping theorem as $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$ is a bounded operator. Note that we always have $\sum_{s'} Z_{s,a,s'} = 0$ for any state-action pair (s, a) , which implies $\tilde{\mathbb{G}}^\pi$ is also in $(M_{\mathcal{F}_{W_1}}^0)^\mathcal{S}$.

Weak Convergence in $\ell^\infty(\mathcal{F}_{KS})$ In the following decomposition

$$\sqrt{n}(\mathcal{I} - \mathcal{T}^\pi)(\hat{\eta}_n^\pi - \eta^\pi) = \underbrace{\sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi}_{(1)} + \underbrace{\sqrt{n} \left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi)}_{(2)},$$

we also have the two relevant terms in $(M_{\mathcal{F}_{KS}}^0)^\mathcal{S}$ under Assumption 1.

Lemma 6.10. Let Assumption 1 hold. For any $s \in \mathcal{S}$, $\sqrt{n} \left[\left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s)$ converge weakly to the process $f \mapsto \tilde{\mathbb{G}}^\pi(s)f$ in $\ell^\infty(\mathcal{F}_{KS})$.

Lemma 6.11. Let Assumption 1 hold. For any $s \in \mathcal{S}$, we have $\left\| \sqrt{n} \left[\left(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) (\hat{\eta}_n^\pi - \eta^\pi) \right] (s) \right\|_{\mathcal{F}_{KS}} = o_P(1)$.

The proof strategy is similar to that of Lemma 6.8 and Lemma 6.9. One may refer to Appendix D for detailed proof.

Because $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$ is not bounded on the space $(M_{\mathcal{F}_{KS}}^0)^\mathcal{S}$, we may not directly apply continuous mapping theorem here. However, the limiting random element indeed lies in a finite-dimensional subspace of $(M_{\mathcal{F}_{KS}}^0)^\mathcal{S}$. To ensure technical rigor, we define $\nu(s, a, s') = \int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a)$, and

$$C_s := \left\{ \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} c_{a,s'} \nu(s, a, s') \mid c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}, \sum_{s' \in \mathcal{S}} c_{a,s'} = 0, \forall a \in \mathcal{A} \right\},$$

and it is straightforward to verify $\tilde{\mathbb{G}}^\pi$ lies in $\prod_{s \in \mathcal{S}} C_s$. Since $\prod_{s \in \mathcal{S}} C_s$ is finite-dimensional, the linear operator $(\mathcal{I} - \mathcal{T}^\pi)^{-1}$ is continuous when confined on $\prod_{s \in \mathcal{S}} C_s$. Finally, we can obtain $\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi) \rightsquigarrow (\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi$.

Weak Convergence in $\ell^\infty(\mathcal{F}_{\text{TV}})$ We again consider the following decomposition

$$\sqrt{n}(\mathcal{I} - \mathcal{T}^\pi)(\hat{\eta}_n^\pi - \eta^\pi) = \underbrace{\sqrt{n}(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)\eta^\pi}_{(1)} + \underbrace{\sqrt{n}(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)(\hat{\eta}_n^\pi - \eta^\pi)}_{(2)}.$$

When Assumption 2 holds, we have the two relevant terms in $(M_{\mathcal{F}_{\text{TV}}}^0)^\mathcal{S}$. The proof idea is nearly identical as before, we first demonstrate term (1) converges weakly to a gaussian random element, then show term (2) is asymptotically negligible. The proofs is in Appendix D.

Lemma 6.12. *Let Assumption 2 hold. For any $s \in \mathcal{S}$, $\sqrt{n}[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)\eta^\pi](s)$ converge weakly to the process $f \mapsto \tilde{\mathbb{G}}^\pi(s)f$ in $\ell^\infty(\mathcal{F}_{\text{TV}})$.*

Lemma 6.13. *Let Assumption 2 hold. For any $s \in \mathcal{S}$, we have $\left\| \sqrt{n}[(\hat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi)(\hat{\eta}_n^\pi - \eta^\pi)](s) \right\|_{\mathcal{F}_{\text{TV}}} = o_P(1)$.*

The final step is also nearly identical to the $\ell^\infty(\mathcal{F}_{\text{KS}})$ case. Since the limiting random element $\tilde{\mathbb{G}}^\pi$ lies in a finite-dimensional subspace of $(M_{\mathcal{F}_{\text{TV}}}^0)^\mathcal{S}$, we can obtain $\sqrt{n}(\hat{\eta}_n^\pi - \eta^\pi) \rightsquigarrow (\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi$ via continuous mapping theorem.

7 Discussions

In this paper, we analyze the statistical performance of distributional reinforcement learning from both non-asymptotic and asymptotic perspectives. Based on our theoretical findings, we devise inferential procedures for a wide class of statistical functionals of the return distribution. We hope our work can spur further research in the uncertainty quantification of reinforcement learning. One future direction is whether we can close the gap between our sample complexity bound $\tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^4}\right)$ and the lower bound $\tilde{O}\left(\frac{1}{\epsilon^2(1-\gamma)^3}\right)$. We speculate that the minimax optimal sample complexity is indeed $O\left(\frac{1}{\epsilon^2(1-\gamma)^3}\right)$, and can be attained through more advanced analysis techniques. Another interesting future direction is to develop non-asymptotic bounds as well as asymptotic results that are uniform for $\pi \in \Pi$, where Π denotes a policy class of interest. This may give rise to a wider range of inferential applications in reinforcement learning.

References

- [1] M. G. Bellemare, W. Dabney, and R. Munos. A distributional perspective on reinforcement learning. In *International conference on machine learning*, pages 449–458. PMLR, 2017.
- [2] M. G. Bellemare, W. Dabney, and M. Rowland. *Distributional Reinforcement Learning*. MIT Press, 2023. <http://www.distributional-rl.org>.
- [3] V. I. Bogachev. *Measure theory*, volume 1. Springer, 2007.
- [4] M. Chae and S. G. Walker. Wasserstein upper bounds of the total variation for smooth densities. *Statistics & Probability Letters*, 163:108771, 2020.
- [5] Y. Chandak, S. Niekum, B. da Silva, E. Learned-Miller, E. Brunskill, and P. S. Thomas. Universal off-policy evaluation. *Advances in Neural Information Processing Systems*, 34: 27475–27490, 2021.
- [6] W. Dabney, G. Ostrovski, D. Silver, and R. Munos. Implicit quantile networks for distributional reinforcement learning. In *International conference on machine learning*, pages 1096–1105. PMLR, 2018.
- [7] W. Dabney, M. Rowland, M. Bellemare, and R. Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [8] W. Dabney, Z. Kurth-Nelson, N. Uchida, C. K. Starkweather, D. Hassabis, R. Munos, and M. Botvinick. A distributional code for value in dopamine-based reinforcement learning. *Nature*, 577(7792):671–675, 2020.
- [9] T. Doan, B. Mazouze, and C. Lyle. Gan q-learning. *arXiv preprint arXiv:1805.04874*, 2018.
- [10] A. Fawzi, M. Balog, A. Huang, T. Hubert, B. Romera-Paredes, M. Barekatin, A. Novikov, F. J. R Ruiz, J. Schrittwieser, G. Swirszcz, et al. Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610(7930):47–53, 2022.
- [11] D. Freirich, T. Shimkin, R. Meir, and A. Tamar. Distributional multivariate policy evaluation and exploration with the bellman gan. In *International Conference on Machine Learning*, pages 1983–1992. PMLR, 2019.

- [12] E. Ghysels, P. Santa-Clara, and R. Valkanov. There is a risk-return trade-off after all. *Journal of financial economics*, 76(3):509–548, 2005.
- [13] B. Hao, X. Ji, Y. Duan, H. Lu, C. Szepesvari, and M. Wang. Bootstrapping fitted q-evaluation for off-policy inference. In *International Conference on Machine Learning*, pages 4074–4084. PMLR, 2021.
- [14] Y. Hua, R. Li, Z. Zhao, X. Chen, and H. Zhang. Gan-powered deep distributional reinforcement learning for resource management in network slicing. *IEEE Journal on Selected Areas in Communications*, 38(2):334–349, 2019.
- [15] A. Huang, L. Leqi, Z. Lipton, and K. Azizzadenesheli. Off-policy risk assessment for markov decision processes. In *International Conference on Artificial Intelligence and Statistics*, pages 5022–5050. PMLR, 2022.
- [16] N. Jiang and L. Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 652–661. PMLR, 2016.
- [17] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- [18] P. W. Lavori and R. Dawson. Dynamic treatment regimes: practical design considerations. *Clinical trials*, 1(1):9–20, 2004.
- [19] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen. Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Advances in neural information processing systems*, 33:12861–12872, 2020.
- [20] X. Li, J. Liang, and Z. Zhang. Online statistical inference for nonlinear stochastic approximation with markovian data. *arXiv preprint arXiv:2302.07690*, 2023.
- [21] X. Li, W. Yang, J. Liang, Z. Zhang, and M. I. Jordan. A statistical analysis of polyak-ruppert averaged q-learning. In *International Conference on Artificial Intelligence and Statistics*, pages 2207–2261. PMLR, 2023.
- [22] S. H. Lim and I. MALIK. Distributional reinforcement learning for risk-sensitive policies. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 30977–30989. Curran

- Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c88a2bd0e793550d0e885aa6e31ca277-Paper-Conference.pdf.
- [23] X. Ma, L. Xia, Z. Zhou, J. Yang, and Q. Zhao. Dsac: Distributional soft actor critic for risk-sensitive reinforcement learning. *arXiv preprint arXiv:2004.14547*, 2020.
- [24] T. Morimura, M. Sugiyama, H. Kashima, H. Hachiya, and T. Tanaka. Nonparametric return distribution approximation for reinforcement learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 799–806, 2010.
- [25] F. Naeem, S. Seifollahi, Z. Zhou, and M. Tariq. A generative adversarial network enabled deep distributional reinforcement learning for transmission scheduling in internet of vehicles. *IEEE Transactions on Intelligent Transportation Systems*, 22(7):4550–4559, 2020.
- [26] OpenAI. Gpt-4 technical report, 2023.
- [27] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [28] N. Ross. Fundamentals of stein’s method. *Probability Surveys*, 8:210–293, 2011.
- [29] M. Rowland, M. Bellemare, W. Dabney, R. Munos, and Y. W. Teh. An analysis of categorical distributional reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 29–37. PMLR, 2018.
- [30] M. Rowland, R. Munos, M. G. Azar, Y. Tang, G. Ostrovski, A. Harutyunyan, K. Tuyls, M. G. Bellemare, and W. Dabney. An analysis of quantile temporal-difference learning. *arXiv preprint arXiv:2301.04462*, 2023.
- [31] C. Shi, S. Zhang, W. Lu, and R. Song. Statistical inference of the value function for reinforcement learning in infinite-horizon settings. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(3):765–793, 2022.
- [32] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.

- [33] H. A. Simon. Dynamic programming under uncertainty with a quadratic criterion function. *Econometrica, Journal of the Econometric Society*, pages 74–81, 1956.
- [34] R. Singh, Q. Zhang, and Y. Chen. Improving robustness via risk averse distributional reinforcement learning. In A. M. Bayen, A. Jadbabaie, G. Pappas, P. A. Parrilo, B. Recht, C. Tomlin, and M. Zeilinger, editors, *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of Machine Learning Research*, pages 958–968. PMLR, 10–11 Jun 2020. URL <https://proceedings.mlr.press/v120/singh20a.html>.
- [35] K. Sun, Y. Zhao, Y. Liu, E. Shi, Y. Wang, X. Yan, B. Jiang, and L. Kong. Interpreting distributional reinforcement learning: A regularization perspective, 2022.
- [36] R. S. Sutton. The reward hypothesis, 2004. URL <http://incompleteideas.net/rlai.cs.ualberta.ca/RLAI/rewardhypothesis.html>.
- [37] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [38] H. Theil. A note on certainty equivalence in dynamic planning. *Econometrica: Journal of the Econometric Society*, pages 346–349, 1957.
- [39] P. Thomas, G. Theodorou, and M. Ghavamzadeh. High-confidence off-policy evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2015.
- [40] A. van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [41] R. Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018. doi: 10.1017/9781108231596.
- [42] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [43] K. Wang, K. Zhou, R. Wu, N. Kallus, and W. Sun. The benefits of being distributional: Small-loss bounds for reinforcement learning. *arXiv preprint arXiv:2305.15703*, 2023.
- [44] R. Wu, M. Uehara, and W. Sun. Distributional offline policy evaluation with predictive error guarantees. *arXiv preprint arXiv:2302.09456*, 2023.

- [45] W. Yang, L. Zhang, and Z. Zhang. Toward theoretical understandings of robust markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):3223–3248, 2022.
- [46] Y. Zhu, J. Dong, and H. Lam. Uncertainty quantification and exploration for reinforcement learning. *Operations Research*, 2023.

A Omitted Proofs in Section 2

Proof of Proposition 2.1. First, we may note that the existence of \mathcal{F} and Q is guaranteed by Kolmogorov's extension theorem. We simply need to verify that $\forall x \in \mathbb{R}$, $\{G^\pi(s) \leq x\}$ is an element of the product σ -algebra \mathcal{F} . Let $G_H^\pi(s) = \sum_{t=0}^H \gamma^t R_t$, then we have

$$\{G^\pi(s) \leq x\} = \bigcap_{H=1}^{\infty} \{G_H^\pi(s) \leq x\}.$$

Since $\forall H$, $\{G_H^\pi(s) \leq x\}$ is an element of the product σ -algebra \mathcal{F} , we also have $\{G^\pi(s) \leq x\} \in \mathcal{F}$. \square

Lemma A.1. *Suppose $\mu \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)^{\mathcal{S}}$ is a vector of distributions and $\mu(s)$ has density $p_s(\cdot)$. If Assumption 1 holds, for any $s \in \mathcal{S}$, $[\mathcal{T}^\pi \mu](s)$ has density $\tilde{p}_s(\cdot)$ such that $\sup_{x \in [0, 1/(1-\gamma)]} \tilde{p}_s(x) \leq C$. Here \mathcal{T}^π can be replaced by any valid distributional Bellman operator.*

Proof. By definition of \mathcal{T}^π , for any $s \in \mathcal{S}$, $[\mathcal{T}^\pi \mu](s)$ also has probability density $\tilde{p}_s(\cdot)$ and

$$\tilde{p}_s(x) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) (p_s(\cdot/\gamma)/\gamma * p_{s,a}^R(\cdot))(x).$$

As for any $x \in [0, 1/(1-\gamma)]$, $|(p_s(\cdot/\gamma)/\gamma * p_{s,a}^R(\cdot))(x)| \leq \sup_{x \in [0, 1/(1-\gamma)]} |p_{s,a}^R(x)|$, we can get $\sup_{x \in [0, 1/(1-\gamma)]} \tilde{p}_s(x) \leq C$ under Assumption 1. \square

Lemma A.2. *If Assumption 1 is true, $\forall s \in \mathcal{S}$, $\eta^\pi(s)$ has density $p_s^G(x)$ being bounded from above by constant C .*

Proof. Define

$$G_H^\pi(s) = \sum_{t=0}^H \gamma^t R_t,$$

where $S_0 = s$, $A_t | S_t \sim \pi(\cdot | S_t)$, $R_t \sim \mathcal{P}_R(S_t, A_t)$ and $S_{t+1} | (S_t, A_t) \sim P(\cdot | S_t, A_t)$. The density of $G_H^\pi(s)$ can be written as

$$p_s^{G_H}(x) = \sum_{\text{all possible length-}H \text{ path } \tau} q(\tau) p_\tau^R(x)$$

with $x \in \left[0, \frac{1}{1-\gamma}\right]$. Suppose $\tau = (s, a_0, s_1, \dots, s_H, a_H)$, then

$$q(\tau) = \pi(a_0 | s) P(s_1 | s, a_0) \pi(a_1 | s_1) P(s_2 | s_1, a_1) \cdots P(s_H | s_{H-1}, a_{H-1}) \pi(a_H | s_H)$$

is the probability of sampling path τ and

$$p_\tau^R(x) = \left[\left(\left(\left(p_{s,a_0}^R(\cdot) * \frac{p_{s_1,a_1}^R(\gamma \cdot)}{\gamma} \right) * \frac{p_{s_2,a_2}^R(\gamma^2 \cdot)}{\gamma^2} \right) \cdots \right) * \frac{p_{s_H,a_H}^R(\gamma^H \cdot)}{\gamma^H} \right] (x)$$

is the density of r.v. $\sum_{t=0}^H \gamma^t (R_t | \tau)$. Here $p_{s,a}^R(x)$ is the density of $\mathcal{P}_R(dr | s, a)$ and $(R_t | \tau) \sim \mathcal{P}_R(\cdot | s_t, a_t)$. According to Lemma E.5, $\sup_x |p_\tau^R(x)| \leq C$. Thus $\sup_x |p_s^{G_H}(x)| \leq C$, *i.e.* $G_H(s)$ has bounded density. This also implies

$$|F_s^{G_H}(x) - F_s^{G_H}(y)| \leq C |x - y|,$$

i.e. the distribution function of $G_H(s)$ is C -Lipschitz continuous. Let $F_s^G(x)$ be the distribution function of $G^\pi(s)$, for any $x \geq y$, we have

$$\begin{aligned} |F_s^G(x) - F_s^G(y)| &= \mathbb{P}(y < G^\pi(s) \leq x) \\ &\leq \mathbb{P}\left(y - \frac{\gamma^H}{1-\gamma} < G_H^\pi(x) \leq x\right) \\ &= \left| F_s^{G_H}(x) - F_s^{G_H}\left(y - \frac{\gamma^H}{1-\gamma}\right) \right| \\ &\leq C \left| x - y + \frac{\gamma^H}{1-\gamma} \right| \\ &= C \left(|x - y| + \frac{\gamma^H}{1-\gamma} \right). \end{aligned}$$

The first inequality is due to the definition of $G_H^\pi(s)$, the second inequality is due to the Lipschitz property of $F_s^{G_H}(x)$. Since H is arbitrarily chosen, we have

$$|F_s^G(x) - F_s^G(y)| \leq C |x - y|,$$

i.e. $G^\pi(s)$ has C -Lipschitz continuous distribution function. As Lipschitz continuity leads to absolute continuity we know $F_s^G(x)$ is absolute continuous. Also, the absolute continuity of $F_s^G(x)$ is equivalent with that the distribution of $G^\pi(s)$ is absolute continuous w.r.t. the Lebesgue measure. By Radon-Nikodym theorem we can get the existence of $p_s^G(x)$. Apparently $\sup_x |p_s^G(x)| \leq C$ by the C -Lipschitz property of $F_s^G(x)$. \square

Proof of Proposition 2.5. For $\eta, \eta' \in \Delta(\mathbb{R})^{\mathcal{S}}$, any $f = \mathbb{1}_A$, $A \in \mathcal{B}(\mathbb{R})$, we have

$$\begin{aligned}
& |[\mathcal{T}^\pi \eta](s)f - [\mathcal{T}^\pi \eta'](s)f| \\
&= \left| \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \pi(a | s) P(s' | s, a) \left[\int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right] \right| \\
&\leq \sup_{a \in \mathcal{A}} \sup_{s' \in \mathcal{S}} \left| \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right| \\
&\leq \sup_{s' \in \mathcal{S}} \sup_{x \in \mathbb{R}} |\eta(s') \mathbb{1}_{(-\infty, x]} - \eta'(s') \mathbb{1}_{(-\infty, x]}| \\
&= \sup_{s' \in \mathcal{S}} \text{KS}(\eta(s'), \eta'(s')).
\end{aligned}$$

The last inequality is due to

$$\begin{aligned}
& \left| \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right| \\
&= \left| \int_0^1 [(b_{r,\gamma})_{\#} \eta(s')] \mathbb{1}_{(-\infty, t]} d\mathcal{P}_R(dr | s, a) - \int_0^1 [(b_{r,\gamma})_{\#} \eta'(s')] \mathbb{1}_{(-\infty, t]} d\mathcal{P}_R(dr | s, a) \right| \\
&= \left| \int_0^1 \eta(s') \mathbb{1}_{(-\infty, (x-r)/\gamma]} d\mathcal{P}_R(dr | s, a) - \int_0^1 \eta'(s') \mathbb{1}_{(-\infty, (x-r)/\gamma]} d\mathcal{P}_R(dr | s, a) \right| \\
&= \left| \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} \eta(s') \mathbb{1}_{(-\infty, (x-R)/\gamma]} - \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} \eta'(s') \mathbb{1}_{(-\infty, (x-R)/\gamma]} \right| \\
&\leq \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} |\eta(s') \mathbb{1}_{(-\infty, (x-R)/\gamma]} - \eta'(s') \mathbb{1}_{(-\infty, (x-R)/\gamma]}| \\
&\leq \sup_{x \in \mathbb{R}} |\eta(s') \mathbb{1}_{(-\infty, x]} - \eta'(s') \mathbb{1}_{(-\infty, x]}|.
\end{aligned}$$

Here $g_{r,\gamma}(A) := \{x \in \mathbb{R} : r + \gamma x \in A\}$ and $g_{r,\gamma}(A) \in \mathcal{B}(\mathbb{R})$ as long as $A \in \mathcal{B}(\mathbb{R})$. Therefore, we have shown that the operator \mathcal{T}^π is non-expansive in the supreme KS distance. Combining Lemma A.1, Lemma A.2 and Proposition 2.2 we have for any $s \in \mathcal{S}$,

$$\text{KS}(\eta^{(k)}(s), \eta^\pi(s)) \leq \sqrt{2CW_1(\eta^{(k)}(s), \eta^\pi(s))}.$$

Combining with Proposition 2.4, we complete the proof. \square

Lemma A.3. *Suppose $\mu \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)^{\mathcal{S}}$ is a vector of distributions and $\mu(s)$ has density $p_s(\cdot) \in H_1^1(\mathbb{R})$. Then under Assumption 2, for any $s \in \mathcal{S}$, $[\mathcal{T}^\pi \mu](s)$ has density $\tilde{p}_s(\cdot) \in H_1^1(\mathbb{R})$. Here \mathcal{T}^π can be replaced by any valid distributional Bellman operator.*

Proof. By definition of \mathcal{T}^π , for any $s \in \mathcal{S}$,

$$\tilde{p}_s(x) = \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) (p_{s,a}^R(\cdot) * p_{s'}(\cdot/\gamma))(x).$$

$$\begin{aligned} \|D^1 \tilde{p}_s\|_1 &\leq \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|D^1 (p_{s,a}^R(\cdot) * p_{s'}(\cdot/\gamma))\|_1 \\ &= \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|(D^1 p_{s,a}^R(\cdot)) * p_{s'}(\cdot/\gamma)\|_1 \\ &\leq \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|D^1 p_{s,a}^R(\cdot)\|_1 \|p_{s'}(\cdot/\gamma)\|_1 \\ &\leq M - 1, \end{aligned}$$

The first equality holds because $(f * g)' = f' * g$, the second inequality holds by Young's convolution inequality (Lemma E.1), and the last inequality holds due to $\|D^1 p_{s,a}^R\|_1 = \|p_{s,a}^R\|_{H_1^1} - \|p_{s,a}^R\|_1 \leq M - 1$. Finally, we may conclude that $\|\tilde{p}_s\|_{H_1^1} = \|D^1 \tilde{p}_s\|_1 + \|\tilde{p}_s\|_1 \leq M$. \square

Lemma A.4. *If Assumption 2 is true, $\forall s \in \mathcal{S}$, $\eta^\pi(s)$ has density $p_s^G(\cdot) \in H_1^1(\mathbb{R})$. Specifically, we have*

$$\|p_s^G\|_{H_1^1} \leq M.$$

Proof. For any $s \in \mathcal{S}$, p_s^G satisfies the following version of distributional Bellman equation:

$$p_s^G(x) = \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) (p_{s,a}^R(\cdot) * p_{s'}^G(\cdot/\gamma))(x).$$

Then we have

$$\begin{aligned} \|D^1 p_s^G\|_1 &\leq \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|D^1 (p_{s,a}^R(\cdot) * p_{s'}^G(\cdot/\gamma))\|_1 \\ &= \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|(D^1 p_{s,a}^R(\cdot)) * p_{s'}^G(\cdot/\gamma)\|_1 \\ &\leq \frac{1}{\gamma} \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) \|D^1 p_{s,a}^R(\cdot)\|_1 \|p_{s'}^G(\cdot/\gamma)\|_1 \\ &\leq M - 1. \end{aligned}$$

The first equality holds because $(f * g)' = f' * g$, the second inequality holds by Young's convolution inequality (Lemma E.1), and the last inequality holds due to $\|D^1 p_{s,a}^R\|_1 = \|p_{s,a}^R\|_{H_1^1} - \|p_{s,a}^R\|_1 \leq M - 1$. Finally, we may conclude that $\|p_s^G\|_{H_1^1} = \|D^1 p_s^G\|_1 + \|p_s^G\|_1 \leq M$. \square

Proof of Proposition 2.7. For $\eta, \eta' \in \Delta(\mathbb{R})^{\mathcal{S}}$, any $f = \mathbb{1}_A$, $A \in \mathcal{B}(\mathbb{R})$, we have

$$\begin{aligned}
& |[\mathcal{T}^\pi \eta](s)f - [\mathcal{T}^\pi \eta'](s)f| \\
&= \left| \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \pi(a | s) P(s' | s, a) \left[\int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right] \right| \\
&\leq \sup_{a \in \mathcal{A}} \sup_{s' \in \mathcal{S}} \left| \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right| \\
&\leq \sup_{s' \in \mathcal{S}} \sup_{A \in \mathcal{B}(\mathbb{R})} |\eta(s') \mathbb{1}_A - \eta'(s') \mathbb{1}_A| \\
&= \sup_{s' \in \mathcal{S}} \text{TV}(\eta(s'), \eta'(s')).
\end{aligned}$$

The last inequality is due to

$$\begin{aligned}
& \left| \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) f - \int_0^1 (b_{r,\gamma})_{\#} \eta'(s') d\mathcal{P}_R(dr | s, a) f \right| \\
&= \left| \int_0^1 [(b_{r,\gamma})_{\#} \eta(s')] \mathbb{1}_A d\mathcal{P}_R(dr | s, a) - \int_0^1 [(b_{r,\gamma})_{\#} \eta'(s')] \mathbb{1}_A d\mathcal{P}_R(dr | s, a) \right| \\
&= \left| \int_0^1 \eta(s') \mathbb{1}_{g_{r,\gamma}(A)} d\mathcal{P}_R(dr | s, a) - \int_0^1 \eta'(s') \mathbb{1}_{g_{r,\gamma}(A)} d\mathcal{P}_R(dr | s, a) \right| \\
&= \left| \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} \eta(s') \mathbb{1}_{g_{r,\gamma}(A)} - \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} \eta'(s') \mathbb{1}_{g_{r,\gamma}(A)} \right| \\
&\leq \mathbb{E}_{R \sim \mathcal{P}_R(\cdot | s, a)} |\eta(s') \mathbb{1}_{g_{r,\gamma}(A)} - \eta'(s') \mathbb{1}_{g_{r,\gamma}(A)}| \\
&\leq \sup_{A \in \mathcal{B}(\mathbb{R})} |\eta(s') \mathbb{1}_A - \eta'(s') \mathbb{1}_A|.
\end{aligned}$$

Therefore, we have shown that the operator \mathcal{T}^π is non-expansive in the supreme TV distance.

For some $\eta \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)$ with Lebesgue density, let p_s^η denote its density function. When Assumption 2 holds true, we have $\|p_s^{\eta^\pi}\|_{H_1^1} \leq M$, $\forall s \in \mathcal{S}$ by Lemma A.4 and $\|p_s^{\eta^{(k)}}\|_{H_1^1} \leq M$ by Lemma A.3. Therefore, for any $s \in \mathcal{S}$,

$$\text{TV}(\eta^{(k)}(s), \eta^\pi(s)) \leq \sqrt{2MKW_1(\eta^{(k)}(s), \eta^\pi(s))}.$$

Combining with Proposition 2.4, we complete the proof. \square

B Omitted Proofs in Section 3

Proof of Corollary 3.1.

$$\max_{s \in \mathcal{S}} W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \leq \left[\frac{1}{(1-\gamma)^{p-1}} \max_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \right]^{\frac{1}{p}}.$$

Therefore we have

$$\begin{aligned} \mathbb{E} \max_{s \in \mathcal{S}} W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s)) &\leq \frac{1}{(1-\gamma)^{1-\frac{1}{p}}} \mathbb{E} \left[\max_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \right]^{\frac{1}{p}} \\ &\leq \frac{1}{(1-\gamma)^{1-\frac{1}{p}}} \left[\mathbb{E} \max_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \right]^{\frac{1}{p}} \\ &\leq \left[\frac{9 \log |\mathcal{S}|}{n(1-\gamma)^{2p+2}} \right]^{\frac{1}{2p}} \end{aligned}$$

and for any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} \max_{s \in \mathcal{S}} W_p(\hat{\eta}_n^\pi(s), \eta^\pi(s)) &\leq \frac{1}{(1-\gamma)^{1-\frac{1}{p}}} \left[\max_{s \in \mathcal{S}} W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \right]^{\frac{1}{p}} \\ &\leq \left[\frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2}}{\sqrt{n(1-\gamma)^{p+1}}} \right]^{\frac{1}{p}}. \end{aligned}$$

□

C Omitted Proofs in Section 4

Proof of Theorem 4.1. Note that $\sqrt{n}W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) = \sup_{f \in \mathcal{F}_{W_1}} \sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s))f$ and $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s)) \rightsquigarrow \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right](s)$ in $\ell^\infty(\mathcal{F}_{W_1})$, thus

$$\sqrt{n}W_1(\hat{\eta}_n^\pi, \eta^\pi) \rightsquigarrow \sup_{f \in \mathcal{F}_{W_1}} \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right](s)f$$

by continuous mapping theorem. It follows that $\lim_{n \rightarrow \infty} \mathbb{P}(W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \in C_1(\alpha)) = 1 - \alpha$.

Also, $\sqrt{n}W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) = \int_0^{\frac{1}{1-\gamma}} |\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s)) \mathbf{1}(-\infty, x]| dx$. and $\sqrt{n}(\hat{\eta}_n^\pi(s) - \eta^\pi(s)) \rightsquigarrow \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right](s)$ in $\ell^\infty(\mathcal{F}_{KS})$ when Assumption 1 holds. Therefore,

$$\sqrt{n}W_1(\hat{\eta}_n^\pi(s), \eta^\pi(s)) \rightsquigarrow \int_0^{\frac{1}{1-\gamma}} \left| \left[(\mathcal{I} - \mathcal{T}^\pi)^{-1} \tilde{\mathbb{G}}^\pi \right](s) \mathbf{1}_{(-\infty, x]} \right| dx$$

by continuous mapping theorem. \square

Proof of Proposition 4.2. Suppose $p \in (\Delta(\mathcal{S}))^{\mathcal{S} \times \mathcal{A}}$ is a transition dynamic, $\eta \in \left(\Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)\right)^{\mathcal{S}}$, and $z \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$. We define $\mathcal{T}^\pi(p)$ as the distributional Bellman operator associated with p , and $\Sigma(p) \in \mathbb{R}^{(\mathcal{S} \times \mathcal{A} \times \mathcal{S}) \times (\mathcal{S} \times \mathcal{A} \times \mathcal{S})}$ as the covariance matrix associated with p defined in Theorem 3.4. Note that the covariance matrix is degenerate such that $\sum_{s,s' \in \mathcal{S}, a \in \mathcal{A}} \left[\Sigma(p)^{\frac{1}{2}} z\right]_{s,a,s'} = 0, \forall z \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$.

$$\begin{aligned} [g_1(\eta, p, z)](s) &:= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} \left[\Sigma(p)^{\frac{1}{2}} z\right]_{s,a,s'} \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) \in M_{\mathcal{F}_{W_1}}^0, \\ g_2(p, \nu) &:= (\mathcal{I} - \mathcal{T}^\pi(p))^{-1} \nu \in \ell^\infty(\mathcal{F}_{W_1}), \\ G(p, \eta, z) &:= \sup_{f \in \mathcal{F}_{W_1}} g_2(p, g_1(\eta, p, z)) f \in \mathbb{R}. \end{aligned}$$

Define the metric $d(\nu_1, \nu_2) := \sup_{s \in \mathcal{S}} \|\nu_1(s) - \nu_2(s)\|_{\mathcal{F}_{W_1}}$ for $\nu_1, \nu_2 \in (\ell^\infty(\mathcal{F}_{W_1}))^{\mathcal{S}}$. Thus given any z , G is continuous in (η, p) when we consider the total variation distance for p and the metric d for ν . According to Lemma E.9, the α -quantile of $G(p, \eta, Z)$ where Z is standard gaussian is continuous in p and η as long as the quantile function of $G(p, \eta, Z)$ is continuous at α . We have 1) z_1 is the quantile function of $G(P, \eta^\pi, Z)$ and \hat{z}_1 is the quantile function of $G(\hat{P}, \hat{\eta}_n^\pi, Z)$; 2) $\hat{P} \xrightarrow{P} P$ and $d(\hat{\eta}_n^\pi, \eta^\pi) \xrightarrow{P} 0$ by Theorem 3.1, the conclusion is true by the continuous mapping theorem. \square

Proof of Proposition 4.4. Let Assumption 1 hold. Suppose $p \in (\Delta(\mathcal{S}))^{\mathcal{S} \times \mathcal{A}}$ is a transition dynamic, $\eta \in \left(\Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)\right)^{\mathcal{S}}$, and $z \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$. We define $\mathcal{T}^\pi(p)$ as the distributional Bellman operator associated with p , and $\Sigma(p) \in \mathbb{R}^{(\mathcal{S} \times \mathcal{A} \times \mathcal{S}) \times (\mathcal{S} \times \mathcal{A} \times \mathcal{S})}$ as the covariance matrix associated with p defined in Theorem 3.4. Note that the covariance matrix is degenerate such that $\sum_{s,s' \in \mathcal{S}, a \in \mathcal{A}} \left[\Sigma(p)^{\frac{1}{2}} z\right]_{s,a,s'} = 0, \forall z \in \mathbb{R}^{\mathcal{S} \times \mathcal{A} \times \mathcal{S}}$.

$$\begin{aligned} [g_1(\eta, p, z)](s) &:= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} \left[\Sigma(p)^{\frac{1}{2}} z\right]_{s,a,s'} \int_0^1 (b_{r,\gamma})_{\#} \eta(s') d\mathcal{P}_R(dr | s, a) \in M_{\mathcal{F}_{\text{KS}}}^0, \\ g_2(p, \nu) &:= (\mathcal{I} - \mathcal{T}^\pi(p))^{-1} \nu \in \ell^\infty(\mathcal{F}_{\text{KS}}), \\ G(p, \eta, z) &:= \sup_{f \in \mathcal{F}_{\text{KS}}} g_2(p, g_1(\eta, p, z)) f \in \mathbb{R}. \end{aligned}$$

Define the metric $d(\nu_1, \nu_2) := \sup_{s \in \mathcal{S}} \|\nu_1(s) - \nu_2(s)\|_{\mathcal{F}_{\text{KS}}}$ for $\nu_1, \nu_2 \in (\ell^\infty(\mathcal{F}_{\text{KS}}))^{\mathcal{S}}$. Thus given any z , G is continuous in (η, p) when we consider the total variation distance for p and the metric d for ν . According to Lemma E.9, the α -quantile of $G(p, \eta, Z)$ where Z is standard gaussian is continuous

in p and η as long as the quantile function of $G(p, \eta, Z)$ is continuous at α . We have 1) z_2 is the quantile function of $G(P, \eta^\pi, Z)$ and \hat{z}_2 is the quantile function of $G(\hat{P}, \hat{\eta}_n^\pi, Z)$; 2) $\hat{P} \xrightarrow{P} P$ and $d(\hat{\eta}_n^\pi, \eta^\pi) \xrightarrow{P} 0$ by Theorem 3.2, the conclusion for z_2 is true by the continuous mapping theorem. The conclusion for z_3 can be shown via similar arguments. \square

D Omitted Proofs in Section 6

Lemma D.1. *Let $F_s(\cdot), \hat{F}_s(\cdot)$ denote the cumulative distribution function of $[\mathcal{T}^\pi \eta^\pi](s), [\hat{\mathcal{T}}_n^\pi \eta^\pi](s)$, respectively. Then $\sqrt{n}(\hat{F}_s(x) - F_s(x))$ is $\frac{1}{\sqrt{n}}$ -sub-gaussian with zero mean, $\forall x \in [0, \frac{1}{1-\gamma}]$, $\forall s \in \mathcal{S}$.*

Proof of Lemma D.1. Recall that

$$\begin{aligned} F_s(x) &= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} w_{s,a,s'} P(s' | s, a) \\ \hat{F}_s(x) &= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} w_{s,a,s'} \hat{P}(s' | s, a), \end{aligned}$$

where the weights $w_{s,a,s'} := \int_0^1 F_{s'}^G\left(\frac{x-r}{\gamma}\right) \mathcal{P}(dr | s, a)$ and we always have $w_{s,a,s'} \in [0, 1]$. Apparently $\mathbb{E}\hat{F}_s(x) = F_s(x)$. By Jensen's inequality, we can get

$$\mathbb{E} \exp\left(\lambda \left(\hat{F}_s(x) - F_s(x)\right)\right) \leq \sum_{a \in \mathcal{A}} \pi(a | s) \mathbb{E} \exp\left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\hat{P}(s' | s, a) - P(s' | s, a)\right)\right).$$

Since

$$\begin{aligned} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\hat{P}(s' | s, a) - P(s' | s, a)\right) &= \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i^{(s,a)} = s'\} - P(s' | s, a)\right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1}\{X_i^{(s,a)} = s'\} - P(s' | s, a)\right). \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{E} \exp\left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\hat{P}(s' | s, a) - P(s' | s, a)\right)\right) \\ &= \mathbb{E} \exp\left(\frac{\lambda}{n} \sum_{i=1}^n \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1}\{X_i^{(s,a)} = s'\} - P(s' | s, a)\right)\right) \\ &= \left[\mathbb{E} \exp\left(\frac{\lambda}{n} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1}\{X_i^{(s,a)} = s'\} - P(s' | s, a)\right)\right)\right]^n. \end{aligned}$$

Because

$$\begin{aligned} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \left\{ X_i^{(s,a)} = s' \right\} - P(s' | s, a) \right) &\leq 1 - \sum_{s' \in \mathcal{S}} w_{s,a,s'} P(s' | s, a) \leq 1, \\ \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \left\{ X_i^{(s,a)} = s' \right\} - P(s' | s, a) \right) &\geq - \sum_{s' \in \mathcal{S}} w_{s,a,s'} P(s' | s, a) \geq -1, \end{aligned}$$

we may obtain

$$\mathbb{E} \exp \left(\frac{\lambda}{n} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \left\{ X_i^{(s,a)} = s' \right\} - P(s' | s, a) \right) \right) \leq \exp \left(\frac{\lambda^2}{2n^2} \right)$$

through Hoeffding's lemma (Lemma E.2). Consequently, we may conclude

$$\mathbb{E} \exp \left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right) \leq \exp \left(\frac{\lambda^2}{2n} \right)$$

and further

$$\mathbb{E} \exp \left(\lambda \left(\widehat{F}_s(x) - F_s(x) \right) \right) \leq \exp \left(\frac{\lambda^2}{2n} \right),$$

which completes the proof. \square

Lemma D.2. *Suppose Assumption 1 holds. Let $p_s(\cdot)$, $\widehat{p}_s(\cdot)$ denote the density function of $[\mathcal{T}^\pi \eta^\pi](s)$, $[\widehat{\mathcal{T}}_n^\pi \eta^\pi](s)$, respectively. Then $\sqrt{n}(\widehat{p}_s(x) - p_s(x))$ is $\frac{C}{\sqrt{n}}$ -sub-gaussian with zero mean, $\forall x \in \left[0, \frac{1}{1-\gamma}\right]$, $\forall s \in \mathcal{S}$.*

Proof of Lemma D.2. Recall that

$$\begin{aligned} p_s(x) &= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} w_{s,a,s'} P(s' | s, a) \\ \widehat{p}_s(x) &= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} w_{s,a,s'} \widehat{P}(s' | s, a), \end{aligned}$$

where the weights $w_{s,a,s'} := \left(p_{s'}^{\gamma G} * p_{s,a}^R \right) (x)$ and we always have $w_{s,a,s'} \in [0, C]$ due to $\sup_{x \in [0, 1/(1-\gamma)]} p_{s,a}^R \leq C$. Apparently $\mathbb{E} \widehat{p}_s(x) = p_s(x)$. By Jensen's inequality, we can get

$$\mathbb{E} \exp \left(\lambda \left(\widehat{p}_s(x) - p_s(x) \right) \right) \leq \sum_{a \in \mathcal{A}} \pi(a | s) \mathbb{E} \exp \left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right).$$

Since

$$\begin{aligned} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) &= \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right). \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{E} \exp \left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right) \\ &= \mathbb{E} \exp \left(\frac{\lambda}{n} \sum_{i=1}^n \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right) \right) \\ &= \left[\mathbb{E} \exp \left(\frac{\lambda}{n} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right) \right) \right]^n. \end{aligned}$$

Because

$$-C \leq \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right) \leq C,$$

we may obtain

$$\mathbb{E} \exp \left(\frac{\lambda}{n} \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\mathbb{1} \{ X_i^{(s,a)} = s' \} - P(s' | s, a) \right) \right) \leq \exp \left(\frac{C^2 \lambda^2}{2n^2} \right)$$

through Hoeffding's lemma (Lemma E.2). Consequently, we may conclude

$$\mathbb{E} \exp \left(\lambda \sum_{s' \in \mathcal{S}} w_{s,a,s'} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right) \leq \exp \left(\frac{C^2 \lambda^2}{2n} \right)$$

and further

$$\mathbb{E} \exp \left(\lambda \left(\widehat{g}_s(x) - g_s(x) \right) \right) \leq \exp \left(\frac{C^2 \lambda^2}{2n} \right),$$

which completes the proof. \square

Proof of Lemma 6.1.

$$\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} = \sup_{s \in \mathcal{S}} \int_0^{\frac{1}{1-\gamma}} \left| \widehat{F}_s(x) - F_s(x) \right| dx.$$

Here $F_s(\cdot)$, $\widehat{F}_s(\cdot)$ denotes the cumulative distribution function of $[\mathcal{T}^\pi \eta^\pi](s)$, $[\widehat{\mathcal{T}}_n^\pi \eta^\pi](s)$, respectively.

Specifically, we have $\widehat{F}_s(x) - F_s(x)$ is $\frac{1}{\sqrt{n}}$ -sub-gaussian (see Lemma D.1) and thus

$$\mathbb{E} \sup_{s \in \mathcal{S}} \left| \widehat{F}_s(x) - F_s(x) \right| \leq 3 \sqrt{\frac{\log |\mathcal{S}|}{n}}, \forall x \in \left[0, \frac{1}{1-\gamma} \right].$$

Therefore, we may get

$$\begin{aligned} \mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} &\leq \int_0^{\frac{1}{1-\gamma}} \mathbb{E} \sup_{s \in \mathcal{S}} \left| \widehat{F}_s(x) - F_s(x) \right| dx \\ &\leq \sqrt{\frac{9 \log |\mathcal{S}|}{n(1-\gamma)^2}}. \end{aligned}$$

The high probability bound for $\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}}$ can be derived via a combination of Lemma 6.1 and McDiarmid's inequality (Lemma E.4). Note that for any fixed $i \in \{1, \dots, n\}$, substituting data vector $\{X_i^{(s,a)}, s \in \mathcal{S}, a \in \mathcal{A}\}$ can change $\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}}$ by at most $\frac{1}{n(1-\gamma)}$. Hence for any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} &\leq \mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{W_1}} + \sqrt{\frac{\log(1/\delta)}{2n(1-\gamma)^2}} \\ &\leq \frac{\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)}/2}{\sqrt{n(1-\gamma)^2}}. \end{aligned}$$

□

Proof of Lemma 6.2. We first verify that the Neumann series $\sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$ converges in $\left(\overline{M_{\mathcal{F}_{W_1}}^0} \right)^{\mathcal{S}}$. First, we claim for any $\nu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^{\mathcal{S}}$, $\left\{ \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu, k = 1, 2, \dots \right\}$ is a Cauchy sequence. WLOG, for any $\nu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^{\mathcal{S}}$, $s \in \mathcal{S}$, we may write $\nu = a_\nu(s) (\nu_+(s) - \nu_-(s))$, where $a_\nu(s)$ is a positive constant and $\nu_+(s), \nu_-(s) \in \Delta \left([0, \frac{1}{1-\gamma}] \right)$. Note that $\|\nu(s)\|_{\mathcal{F}_{W_1}} = a_\nu(s) W_1(\nu_+(s), \nu_-(s))$ for any

$s \in \mathcal{S}$. For $k_1 < k_2$, we have for any $s \in \mathcal{S}$

$$\begin{aligned}
\left\| \sum_{i=k_1}^{k_2} [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} &\leq \sum_{i=k_1}^{k_2} \left\| [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&= \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} W_1 \left([(\mathcal{T}^\pi)^i \nu_+] (s), [(\mathcal{T}^\pi)^i \nu_-] (s) \right) \\
&\leq \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \gamma^i \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s)) \\
&\leq \frac{\gamma^{k_1} \sup_{s \in \mathcal{S}} a_\nu(s) \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s))}{1 - \gamma}.
\end{aligned}$$

The second last inequality is due to Proposition 2.4. Then we define ν_* such that for any $f \in \mathcal{F}_{W_1}$, $\nu_*(s)f := \lim_{k \rightarrow \infty} \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu] (s)f$. ν_* is well-defined because $\left\{ \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu(s)f, k = 1, 2, \dots \right\}$ is a Cauchy sequence in \mathbb{R} by similar arguments as above. As

$$\begin{aligned}
\left\| \nu_*(s) - \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} &= \left\| \sum_{i=k}^{\infty} [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&\leq \sum_{i=k}^{\infty} \left\| [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&\leq \frac{\gamma^k \sup_{s \in \mathcal{S}} a_\nu(s) \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s))}{1 - \gamma},
\end{aligned}$$

$\nu_* = \lim_{k \rightarrow \infty} \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu$. Apparently $\nu_* \in \left(\overline{M_{\mathcal{F}_{W_1}}^0} \right)^\mathcal{S}$. Since for any $\nu \in \left(M_{\mathcal{F}_{W_1}}^0 \right)^\mathcal{S}$,

$$\begin{aligned}
&\left\| [(\mathcal{I} - \mathcal{T}^\pi)^{-1} \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&= \left\| \sum_{i=0}^{\infty} [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&\leq \sum_{i=0}^{\infty} \left\| [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{W_1}} \\
&\leq \sum_{i=0}^{\infty} \gamma^i \sup_{s \in \mathcal{S}} a_\nu(s) W_1 (\nu_+(s), \nu_-(s)) \\
&\leq \frac{\sup_{s \in \mathcal{S}} a_\nu(s) W_1 (\nu_+(s), \nu_-(s))}{1 - \gamma} \\
&= \frac{\sup_{s \in \mathcal{S}} \|\nu(s)\|_{\mathcal{F}_{W_1}}}{1 - \gamma}.
\end{aligned}$$

We complete the proof. We comment that proof is inspired by the B.L.T. theorem (continuous linear extension theorem) in functional analysis. \square

Proof of Lemma 6.3. When Assumption 1 holds,

$$\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} = \sup_{s \in \mathcal{S}} \frac{1}{2} \int_0^{\frac{1}{1-\gamma}} |\hat{p}_s(x) - p_s(x)| dx.$$

Here $p_s(\cdot)$, $\hat{p}_s(\cdot)$ denotes the density function of $[\mathcal{T}^\pi \eta^\pi](s)$, $[\widehat{\mathcal{T}}_n^\pi \eta^\pi](s)$, respectively. Specifically, we have $\hat{p}_s(x) - p_s(x)$ is $\frac{C}{\sqrt{n}}$ -sub-gaussian (see Lemma D.2) and thus

$$\mathbb{E} \sup_{s \in \mathcal{S}} |\hat{p}_s(x) - p_s(x)| \leq 3C \sqrt{\frac{\log |\mathcal{S}|}{n}}, \forall x \in \left[0, \frac{1}{1-\gamma} \right].$$

Therefore, we may get

$$\begin{aligned} \mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} &\leq \frac{1}{2} \int_0^{\frac{1}{1-\gamma}} \mathbb{E} \sup_{s \in \mathcal{S}} |\hat{p}_s(x) - p_s(x)| dx \\ &\leq \frac{3C}{2} \sqrt{\frac{\log |\mathcal{S}|}{n(1-\gamma)^2}}. \end{aligned}$$

The high probability bound for $\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}$ can be derived via a combination of Lemma 6.1 and McDiarmid's inequality (Lemma E.4). Note that for any fixed $i \in \{1, \dots, n\}$, substituting data vector $\{X_i^{(s,a)}, s \in \mathcal{S}, a \in \mathcal{A}\}$ can change $\sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}$ by at most $\frac{C}{2n(1-\gamma)}$. Hence for any $\delta > 0$, with probability at least $1 - \delta$,

$$\begin{aligned} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} &\leq \mathbb{E} \sup_{s \in \mathcal{S}} \left\| \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} + \frac{C}{2} \sqrt{\frac{\log(1/\delta)}{2n(1-\gamma)^2}} \\ &\leq \frac{C \left(\sqrt{9 \log |\mathcal{S}|} + \sqrt{\log(1/\delta)/2} \right)}{2\sqrt{n(1-\gamma)^2}}. \end{aligned}$$

\square

Proof of Lemma 6.4. It suffices to verify that the Neumann series $\sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$ converges in $\left(\overline{M_{\mathcal{F}_{\text{KS}}}^0} \right)^\mathcal{S}$. First, we claim for any $\nu \in \left(M_{\mathcal{F}_{\text{KS}}}^0 \right)^\mathcal{S}$, $\left\{ \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu, k = 1, 2, \dots \right\}$ is a Cauchy sequence. WLOG, for any $\nu \in \left(M_{\mathcal{F}_{\text{KS}}}^0 \right)^\mathcal{S}$, $s \in \mathcal{S}$, we may write $\nu = a_\nu(s) (\nu_+(s) - \nu_-(s))$, where $a_\nu(s)$ is a positive constant and $\nu_+(s), \nu_-(s) \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right)$. Note that $\|\nu(s)\|_{\mathcal{F}_{\text{KS}}} = a_\nu(s) \text{KS}(\nu_+(s), \nu_-(s))$ for any

$s \in \mathcal{S}$. There also exists a positive constant C_ν such that suppose for any $s \in \mathcal{S}$, both $\nu_+(s)$ and $\nu_-(s)$ have a density bounded by C_ν . For $k_1 < k_2$, we have for any $s \in \mathcal{S}$

$$\begin{aligned}
\left\| \sum_{i=k_1}^{k_2} [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{\text{KS}}} &\leq \sum_{i=k_1}^{k_2} \left\| [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&= \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \text{KS} \left([(\mathcal{T}^\pi)^i \nu_+] (s), [(\mathcal{T}^\pi)^i \nu_-] (s) \right) \\
&\leq \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \sqrt{2C_\nu W_1 \left((\mathcal{T}^\pi)^i \nu_+(s), (\mathcal{T}^\pi)^i \nu_-(s) \right)} \\
&\leq \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \sqrt{2C_\nu \gamma^i \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s))} \\
&\leq \frac{\gamma^{k_1/2} \sup_{s \in \mathcal{S}} a_\nu(s) \sqrt{2C_\nu \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s))}}{1 - \sqrt{\gamma}}.
\end{aligned}$$

The second inequality is by Proposition 2.2 and the third inequality is due to Proposition 2.4. Then we define ν_* such that for any $f \in \mathcal{F}_{\text{KS}}$, $\nu_*(s)f := \lim_{k \rightarrow \infty} \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu] (s)f$. ν_* is well-defined because $\left\{ \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu(s)f, k = 1, 2, \dots \right\}$ is a Cauchy sequence in \mathbb{R} by similar arguments as above.

As

$$\begin{aligned}
\left\| \nu_*(s) - \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{\text{KS}}} &= \left\| \sum_{i=k}^{\infty} [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&\leq \sum_{i=k}^{\infty} \left\| [(\mathcal{T}^\pi)^i \nu] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\
&\leq \frac{\gamma^{k/2} \sup_{s \in \mathcal{S}} a_\nu(s) \sqrt{2C_\nu \sup_{s \in \mathcal{S}} W_1 (\nu_+(s), \nu_-(s))}}{1 - \sqrt{\gamma}},
\end{aligned}$$

$\nu_* = \lim_{k \rightarrow \infty} \sum_{i=1}^k (\mathcal{T}^\pi)^i \nu$. Apparently $\nu_* \in \left(\overline{M_{\mathcal{F}_{\text{KS}}}^0} \right)^\mathcal{S}$. \square

Proof of Proposition 6.2. For any $x \in \left[0, \frac{1}{1-\gamma} \right]$, we define $h_\epsilon^x(\cdot)$ to be a smoothed version of $\mathbf{1}_{(-\infty, x]}$. Specifically,

$$h_\epsilon^x(t) = \begin{cases} 1 & t \leq x \\ 1 - \frac{t-x}{\epsilon} & x < t \leq x + \epsilon \\ 0 & t > x + \epsilon \end{cases}$$

We have $h_\epsilon^x(t)$ is Lipschitz with Lipschitz coefficient $\frac{1}{\epsilon}$ and $h_\epsilon^x(t) \geq \mathbf{1}_{(-\infty, x]}$. For any $\mu \in M_{\mathcal{F}_{\text{KS}}}^0$, we

have the Jordan decomposition $\mu = \mu_+ - \mu_-$. Then we have

$$\begin{aligned}
\mu \mathbf{1}_{(-\infty, x]} &= \mu_+ \mathbf{1}_{(-\infty, x]} - \mu_- \mathbf{1}_{(-\infty, x]} \\
&= \mu_+ \mathbf{1}_{(-\infty, x]} - \mu_- h_\epsilon^x(\cdot) + \mu_- h_\epsilon^x(\cdot) - \mu_- \mathbf{1}_{(-\infty, x]} \\
&\leq (\mu_+ - \mu_-) h_\epsilon^x(\cdot) + \mu_- (h_\epsilon^x(\cdot) - \mathbf{1}_{(\infty, x]}) \\
&\leq \frac{1}{\epsilon} \mu(\epsilon h_\epsilon^x(\cdot)) + \mu_-(h_\epsilon^x(\cdot) - \mathbf{1}_{(\infty, x]}) \\
&\leq \frac{\|\mu\|_{\mathcal{F}_{W_1}}}{\epsilon} + \epsilon \|\mu\|_{\text{loc}}.
\end{aligned}$$

Setting $\epsilon = \frac{\|\mu\|_{\mathcal{F}_{W_1}}}{\|\mu\|_{\text{loc}}}$ yields

$$\mu \mathbf{1}_{(-\infty, x]} \leq \sqrt{2 \|\mu\|_{\text{loc}} \|\mu\|_{\mathcal{F}_{W_1}}}.$$

We may also prove $-\mu \mathbf{1}_{(-\infty, x]} \leq \sqrt{2 \|\mu\|_{\text{loc}} \|\mu\|_{\mathcal{F}_{W_1}}}$ by similar arguments. \square

Proof of Lemma 6.5. We write the Jordan decomposition of $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$ as

$$\left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) = a(s) (\mu_+(s) - \mu_-(s))$$

Here for any $s \in \mathcal{S}$, $\mu_+(s), \mu_-(s) \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right)$ have probability density functions $p_+(s), p_-(s)$. And it is easy to verify $a(s) = \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\mathcal{F}_{\text{TV}}}$. Recall that for any $\nu \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right)^\mathcal{S}$ with $p_s^\nu(\cdot)$ as the probability density function of $\nu(s)$, we have for any valid distributional Bellman operator \mathcal{T}^π , $[\mathcal{T}^\pi \nu] (s)$ also has probability density $\tilde{p}_s(\cdot)$ and

$$\tilde{p}_s(x) = \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} P(s' | s, a) (p_{s'}^\nu(\cdot/\gamma)/\gamma * p_{s,a}^R(\cdot)) (x).$$

As for any $x \in \left[0, \frac{1}{1-\gamma} \right]$, $|(p_{s'}^\nu(\cdot/\gamma)/\gamma * p_{s,a}^R(\cdot)) (x)| \leq \sup_{x \in [0, 1/(1-\gamma)]} |p_{s,a}^R(x)|$, we can get $\|[\mathcal{T}^\pi \nu] (s)\|_{\text{loc}} \leq C$ under Assumption 1. Therefore, for $i \geq 1$, we have

$$\begin{aligned}
\left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right] (s) \right\|_{\text{loc}} &\leq \sup_{s \in \mathcal{S}} a(s) \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\mu_+ - \mu_-) \right] (s) \right\|_{\text{loc}} \\
&\leq C \sup_{s \in \mathcal{S}} a(s)
\end{aligned}$$

We complete the proof. \square

Proof of Proposition 6.3. For any $\mu \in \left(M_{\mathcal{F}_{\text{TV}}}^0\right)^{\mathcal{S}}$, we may write $\mu(s) = a_\mu(s)(\mu_+(s) - \mu_-(s))$, where $\mu_+(s), \mu_-(s) \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)$ with density $p_s^+(\cdot), p_s^-(\cdot) \in H_1^1(\mathbb{R})$, and $a_\mu(s)$ is a normalization factor. Let $\{\phi_m\}$ be the orthogonal system in $L^2([-1, 1])$ of Legendre polynomials, and kernel $K_h(x) = h^{-1}(\phi_0(0)\phi_0(x/h) + \phi_1(0)\phi_1(x/h))$. From Lemma 2.1 in [4],

$$\begin{aligned}
\|\mu(s)\|_{\mathcal{F}_{\text{TV}}} &= \frac{1}{2}a_\mu(s) \|p_s^+ - p_s^-\|_1 \\
&= \frac{1}{2}a_\mu(s) (\|p_s^+ - K_h * p_s^+\|_1 + \|K_h * p_s^- - K_h * p_s^+\|_1 + \|p_s^- - K_h * p_s^-\|_1) \\
&\leq K'a_\mu(s)(h \|p_s^+\|_{H_1^1} + h \|p_s^-\|_{H_1^1} + W_1(p_s^+, p_s^-)/h) \\
&= K'h \left[a_\mu(s) (\|p_s^+\|_{H_1^1} + \|p_s^-\|_{H_1^1}) \right] + \frac{K'}{h} a_\mu(s) W_1(p_s^+, p_s^-) \\
&= K'h \|\mu(s)\|_{H_1^1} + \frac{K'}{h} \|\mu(s)\|_{\mathcal{F}_{W_1}}.
\end{aligned}$$

The conclusion follows by choosing $h = \sqrt{\|\mu(s)\|_{\mathcal{F}_{W_1}} / (2K' \|\mu(s)\|_{H_1^1})}$. \square

Lemma 6.6. It suffices to verify that the Neumann series $\sum_{i=0}^{\infty} (\mathcal{T}^\pi)^i$ converges in $\left(\overline{M_{\mathcal{F}_{\text{TV}}}^0}\right)^{\mathcal{S}}$. First, we claim for any $\nu \in \left(M_{\mathcal{F}_{\text{TV}}}^0\right)^{\mathcal{S}}$, $\left\{\sum_{i=1}^k (\mathcal{T}^\pi)^i \nu, k = 1, 2, \dots\right\}$ is a Cauchy sequence. WLOG, for any $\nu \in \left(M_{\mathcal{F}_{\text{TV}}}^0\right)^{\mathcal{S}}$, $s \in \mathcal{S}$, we may write $\nu = a_\nu(s)(\nu_+(s) - \nu_-(s))$, where $a_\nu(s)$ is a positive constant and $\nu_+(s), \nu_-(s) \in \Delta\left(\left[0, \frac{1}{1-\gamma}\right]\right)$. Note that $\|\nu(s)\|_{\mathcal{F}_{\text{TV}}} = a_\nu(s)\text{TV}(\nu_+(s), \nu_-(s))$ for any $s \in \mathcal{S}$. There also exists a positive constant M_ν such that for any $s \in \mathcal{S}$, both $\nu_+(s)$ and $\nu_-(s)$ have a density in $H_1^1(\mathbb{R})$ and $\|\nu_+(s)\|_{H_1^1} \leq M_\nu$, $\|\nu_-(s)\|_{H_1^1} \leq M_\nu$. For $k_1 < k_2$, we have for any $s \in \mathcal{S}$

$$\begin{aligned}
\left\| \sum_{i=k_1}^{k_2} [(\mathcal{T}^\pi)^i \nu](s) \right\|_{\mathcal{F}_{\text{TV}}} &\leq \sum_{i=k_1}^{k_2} \left\| [(\mathcal{T}^\pi)^i \nu](s) \right\|_{\mathcal{F}_{\text{TV}}} \\
&= \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \text{TV} \left([(\mathcal{T}^\pi)^i \nu_+](s), [(\mathcal{T}^\pi)^i \nu_-](s) \right) \\
&\leq \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \sqrt{2K M_\nu W_1 \left((\mathcal{T}^\pi)^i \nu_+(s), (\mathcal{T}^\pi)^i \nu_-(s) \right)} \\
&\leq \sup_{s \in \mathcal{S}} a_\nu(s) \sum_{i=k_1}^{k_2} \sqrt{2K M_\nu \gamma^i \sup_{s \in \mathcal{S}} W_1(\nu_+(s), \nu_-(s))} \\
&\leq \frac{\gamma^{k_1/2} \sup_{s \in \mathcal{S}} a_\nu(s) \sqrt{2K M_\nu \sup_{s \in \mathcal{S}} W_1(\nu_+(s), \nu_-(s))}}{1 - \sqrt{\gamma}}.
\end{aligned}$$

The second inequality is by Proposition 2.3 and the third inequality is due to Proposition 2.4. Then

we define ν_* such that for any $f \in \mathcal{F}_{\text{TV}}$, $\nu_*(s)f := \lim_{k \rightarrow \infty} \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu](s)f$. ν_* is well-defined because $\left\{ \sum_{i=0}^k (\mathcal{T}^\pi)^i \nu(s)f, k = 1, 2, \dots \right\}$ is a Cauchy sequence in \mathbb{R} by similar arguments as above.

As

$$\begin{aligned} \left\| \nu_*(s) - \sum_{i=0}^k [(\mathcal{T}^\pi)^i \nu](s) \right\|_{\mathcal{F}_{\text{TV}}} &= \left\| \sum_{i=k}^{\infty} [(\mathcal{T}^\pi)^i \nu](s) \right\|_{\mathcal{F}_{\text{TV}}} \\ &\leq \sum_{i=k}^{\infty} \left\| [(\mathcal{T}^\pi)^i \nu](s) \right\|_{\mathcal{F}_{\text{TV}}} \\ &\leq \frac{\gamma^{k/2} \sup_{s \in \mathcal{S}} a_\nu(s) \sqrt{2KM_\nu \sup_{s \in \mathcal{S}} W_1(\nu_+(s), \nu_-(s))}}{1 - \sqrt{\gamma}}, \end{aligned}$$

$\nu_* = \lim_{k \rightarrow \infty} \sum_{i=0}^k (\mathcal{T}^\pi)^i \nu$. Apparently $\nu_* \in \left(\overline{M_{\mathcal{F}_{\text{TV}}}^0} \right)^\mathcal{S}$. \square

Proof of Lemma 6.7. We write the Jordan decomposition of $(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi$ as

$$\left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right](s) = a(s)(\mu_+(s) - \mu_-(s))$$

Here for any $s \in \mathcal{S}$, $\mu_+(s), \mu_-(s) \in \Delta \left(\left[0, \frac{1}{1-\gamma} \right] \right)$ has probability density function $p_+(s), p_-(s) \in H_1^1(\mathbb{R})$. And it is easy to verify $a(s) = \left\| \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right](s) \right\|_{\mathcal{F}_{\text{TV}}}$. For $i \geq 1$, we have

$$\begin{aligned} \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right](s) \right\|_{H_1^1} &\leq \sup_{s \in \mathcal{S}} a(s) \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i (\mu_+ - \mu_-) \right](s) \right\|_{H_1^1} \\ &= \sup_{s \in \mathcal{S}} a(s) \left(\left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i \mu_+ \right](s) \right\|_{H_1^1} + \left\| \left[(\widehat{\mathcal{T}}_n^\pi)^i \mu_- \right](s) \right\|_{H_1^1} \right) \\ &\leq 2 \sup_{s \in \mathcal{S}} a(s) M. \end{aligned}$$

The last inequality holds by Lemma A.3. \square

Proof of Lemma 6.8. For any $s \in \mathcal{S}$, we have

$$\begin{aligned} \sqrt{n} \left[(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi) \eta^\pi \right](s) &= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \int_0^1 (b_{r,\gamma})_\# \eta^\pi(s') \mathcal{P}_R(dr | s, a) \\ &= \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 (b_{r,\gamma})_\# \eta^\pi(s') \mathcal{P}_R(dr | s, a) \right] \end{aligned}$$

Thus the conclusion follows via the multivariate CLT and Lemma E.7. \square

Proof of Lemma 6.9. For any $s \in \mathcal{S}$, we have

$$\begin{aligned}
& \left\| \sqrt{n} \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \left(\widehat{\eta}_n^\pi - \eta^\pi \right) \right] (s) \right\|_{\mathcal{F}_{W_1}} \\
&= \left\| \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right] \right\|_{\mathcal{F}_{W_1}} \\
&\leq \left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1 \sup_{a \in \mathcal{A}, s' \in \mathcal{S}} \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{W_1}}.
\end{aligned}$$

Since $\left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1$ is of the order $O_P(1)$, it suffices to show for any $s' \in \mathcal{S}, a \in \mathcal{A}$, $\left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{W_1}} = o_P(1)$. Noting that

$$\begin{aligned}
& \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{W_1}} \\
&= W_1 \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right).
\end{aligned}$$

We claim that

$$W_1 \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right) \leq W_1 \left(\widehat{\eta}_n^\pi(s'), \eta^\pi(s') \right).$$

For simplicity of notations, we use $\widehat{\nu}$ in short of $\int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a)$ and ν in short of $\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a)$. In fact, suppose two random variables $X \sim \eta^\pi(s')$, $Y \sim \widehat{\eta}_n^\pi(s')$, and an independent random variable $R \sim \mathcal{P}_R(dr | s, a)$, then $R + \gamma Y \sim \widehat{\nu}$ and $R + \gamma X \sim \nu$. Then we have

$$\begin{aligned}
& W_1 \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right) \\
&= \inf_{W \sim \nu, Z \sim \widehat{\nu}} \mathbb{E} |W - Z| \\
&\leq \mathbb{E} |(R + \gamma X) - (R + \gamma Y)| \\
&= \gamma \mathbb{E} |X - Y|.
\end{aligned}$$

Our claim is true since X and Y are chosen arbitrarily. Our conclusion follows since by Theorem 3.1

$$W_1(\eta^\pi(s'), \widehat{\eta}_n^\pi(s')) = o_P(1). \quad \square$$

Proof of Lemma 6.10. The proof is identical to that of Lemma 6.8. For any $s \in \mathcal{S}$, we have

$$\sqrt{n} \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) = \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a) \right]$$

Thus the conclusion follows via the multivariate CLT and Lemma E.7. \square

Proof of Lemma 6.11. For any $s \in \mathcal{S}$, we have

$$\begin{aligned} & \left\| \sqrt{n} \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \left(\widehat{\eta}_n^\pi - \eta^\pi \right) \right] (s) \right\|_{\mathcal{F}_{\text{KS}}} \\ &= \left\| \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right] \right\|_{\mathcal{F}_{\text{KS}}} \\ &\leq \left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1 \sup_{a \in \mathcal{A}, s' \in \mathcal{S}} \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{KS}}}. \end{aligned}$$

Since $\left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1$ is of the order $O_P(1)$, it suffices to show for any $s' \in \mathcal{S}, a \in \mathcal{A}$, $\left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{KS}}} = o_P(1)$. Noting that

$$\begin{aligned} & \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{KS}}} \\ &= \text{KS} \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right). \end{aligned}$$

We claim that

$$\text{KS} \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right) \leq \text{KS} \left(\widehat{\eta}_n^\pi(s'), \eta^\pi(s') \right).$$

The claim can be verified using the same argument as in the proof of Proposition 2.5 where we show the operator \mathcal{T}^π is non-expansive in supreme KSdistance. Our conclusion follows since by Theorem 3.2 $\text{KS}(\eta^\pi(s'), \widehat{\eta}_n^\pi(s')) = o_P(1)$. \square

Proof of Lemma 6.12. The proof is identical to that of Lemma 6.8. For any $s \in \mathcal{S}$, we have

$$\sqrt{n} \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \eta^\pi \right] (s) = \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a) \right]$$

Thus the conclusion follows via the multivariate CLT and Lemma E.7. \square

Proof of Lemma 6.13. For any $s \in \mathcal{S}$, we have

$$\begin{aligned}
& \left\| \sqrt{n} \left[\left(\widehat{\mathcal{T}}_n^\pi - \mathcal{T}^\pi \right) \left(\widehat{\eta}_n^\pi - \eta^\pi \right) \right] (s) \right\|_{\mathcal{F}_{\text{TV}}} \\
&= \left\| \sum_{a \in \mathcal{A}, s' \in \mathcal{S}} \left[\sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right] \left[\pi(a | s) \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right] \right\|_{\mathcal{F}_{\text{TV}}} \\
&\leq \left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1 \sup_{a \in \mathcal{A}, s' \in \mathcal{S}} \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{TV}}}.
\end{aligned}$$

Since $\left\| \sqrt{n} \left(\widehat{P}(s' | s, a) - P(s' | s, a) \right) \right\|_1$ is of the order $O_P(1)$, it suffices to show for any $s' \in \mathcal{S}, a \in \mathcal{A}$, $\left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{TV}}} = o_P(1)$. Noting that

$$\begin{aligned}
& \left\| \int_0^1 \left[(b_{r,\gamma})_{\#} \eta^\pi(s') - (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \right] \mathcal{P}_R(dr | s, a) \right\|_{\mathcal{F}_{\text{TV}}} \\
&= \text{TV} \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right).
\end{aligned}$$

We claim that

$$\text{TV} \left(\int_0^1 (b_{r,\gamma})_{\#} \eta^\pi(s') \mathcal{P}_R(dr | s, a), \int_0^1 (b_{r,\gamma})_{\#} \widehat{\eta}_n^\pi(s') \mathcal{P}_R(dr | s, a) \right) \leq \text{TV}(\widehat{\eta}_n^\pi(s'), \eta^\pi(s')).$$

The claim can be verified using the same argument as in the proof of Proposition 2.7 where we show the operator \mathcal{T}^π is non-expansive in supreme TVdistance. Our conclusion follows by Theorem 3.3 $\text{TV}(\eta^\pi(s'), \widehat{\eta}_n^\pi(s')) = o_P(1)$. \square

E Technical Lemmas

Lemma E.1 (Young's Convolution Inequality). *Suppose f is in the Lebesgue space $L^p(\mathbb{R}^d)$ and g is in $L^q(\mathbb{R}^d)$ and $\frac{1}{p} + \frac{1}{q} = \frac{1}{r} + 1$ with $1 \leq p, q, r \leq \infty$. Then $\|f\|_p \|g\|_q \leq \|f * g\|_r$.*

Proof. See Theorem 3.9.4 in [3]. \square

Lemma E.2 (Hoeffding's Lemma). *Suppose $X \in [a, b]$ is a random variable with $\mathbb{E}X = 0$, then for any $\lambda \in \mathbb{R}$,*

$$\mathbb{E} \exp(\lambda X) \leq \exp\left(\frac{\lambda^2(b-a)^2}{8}\right)$$

Proof. By Jensen's inequality, for any $x \in [a, b]$,

$$\exp(\lambda x) \leq \frac{b-x}{b-a} \exp(\lambda a) + \frac{x-a}{b-a} \exp(\lambda b).$$

So

$$\mathbb{E} \exp(\lambda X) \leq \frac{b}{b-a} \exp(\lambda a) - \frac{a}{b-a} \exp(\lambda b) := \exp(L(\lambda(b-a))),$$

where $L(h) = \frac{ha}{b-a} + \log\left(1 + \frac{(1-e^h)a}{b-a}\right)$. We may find $L(0) = L'(0) = 0$ and $L''(h) = -\frac{abe^h}{(b-ae^h)^2} \leq \frac{1}{4}$.

By Taylor's theorem, there is some $\theta \in [0, 1]$ such that

$$L(h) = L(0) + L'(0)h + \frac{1}{2}h^2L''(\theta h) \leq \frac{1}{8}h^2.$$

We complete the proof. □

Lemma E.3 (Hoeffding's Inequality). *Suppose $\{X_1, \dots, X_n\}$ be n i.i.d. random variables with values in $[0, 1]$. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ we have*

$$\left| \mathbb{E}X_1 - \frac{\sum_{i=1}^n X_i}{n} \right| \leq \sqrt{\frac{\log(2/\delta)}{2n}}.$$

Proof. See Theorem 2.2.6 in [41]. □

Lemma E.4 (McDiarmid's inequality). *Let $f: \mathcal{X}_1 \times \dots \times \mathcal{X}_N$ satisfy the bounded differences property with bounds c_1, \dots, c_N , i.e. there are constants c_1, \dots, c_N such that for all $i \in \{1, \dots, N\}$, $x_1 \in \mathcal{X}_1, \dots, x_N \in \mathcal{X}_N$, $\sup_{x'_i \in \mathcal{X}_i} |f(x_1, \dots, x_i, \dots, x_N) - f(x_1, \dots, x'_i, \dots, x_N)| \leq c_i$. Consider independent random variables X_1, \dots, X_N where $X_i \in \mathcal{X}_i$ for all i , then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,*

$$|f(X_1, \dots, X_N) - \mathbb{E}f(X_1, \dots, X_N)| \leq \sqrt{\frac{\left(\sum_{i=1}^N c_i^2\right) \log(2/\delta)}{2}}.$$

Proof. See Theorem 2.9.1 in [41]. □

Lemma E.5. *Suppose X_1, \dots, X_N are a sequence of independent r.v.s, X_i has density $p_i(x)$ and $\sum_{i=1}^N X_i$ has density $p(x)$. If $\sup_{x \in \mathbb{R}} p_1(x) \leq C$, then $\sup_{x \in \mathbb{R}} p(x) \leq C$.*

Proof. We have

$$p(x) = [(((p_1 * p_2) * p_3) \cdots) * p_N](x),$$

where $[f * g](x) := \int_{\mathbb{R}} f(x-y)g(y)dy$. Therefore, $[p_1 * p_2](x) = \mathbb{E}p_1(x-X_2)$ and $\sup_{x \in \mathbb{R}} [p_1 * p_2](x) \leq C$. We may complete the proof by deduction. \square

Lemma E.6. *Let X_1, \dots, X_N be any $N \geq 2$ random variables such that for any $\lambda > 0$,*

$$\mathbb{E} \exp(\lambda X_i) \leq \exp\left(\frac{\sigma^2 \lambda^2}{2}\right).$$

Then

$$\mathbb{E} \max_{i \in \{1, \dots, N\}} |X_i| \leq \sqrt{9\sigma^2 \log N}.$$

Proof. Fix $t > 0$, for any $\lambda > 0$,

$$\begin{aligned} \mathbb{P}(X_i > t) &\leq \frac{\mathbb{E} \exp(\lambda X_i)}{\exp(\lambda t)} \\ &\leq \exp\left(\frac{\sigma^2 \lambda^2}{2} - \lambda t\right). \end{aligned}$$

Setting $\lambda = \frac{t}{\sigma^2}$ yields

$$\mathbb{P}(X_i > t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)$$

and

$$\mathbb{P}(|X_i| > t) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

Therefore,

$$\mathbb{P}\left(\frac{|X_i|}{\sqrt{1 + \log i}} > t\right) \leq 2 \exp\left(-\frac{(1 + \log i)t^2}{2\sigma^2}\right)$$

and

$$\mathbb{P}\left(\max_{i \in \{1, \dots, N\}} \frac{|X_i|}{\sqrt{1 + \log i}} > t\right) \leq \sum_{i=1}^N 2 \exp\left(-\frac{(1 + \log i)t^2}{2\sigma^2}\right) = 2 \exp\left(-\frac{t^2}{2\sigma^2}\right) \sum_{i=1}^N i^{-\frac{t^2}{2\sigma^2}}$$

Therefore,

$$\begin{aligned} \mathbb{E} \max_{i \in \{1, \dots, N\}} \frac{|X_i|}{\sqrt{1 + \log i}} &= \int_0^\infty \mathbb{P}\left(\max_{i \in \{1, \dots, N\}} \frac{|X_i|}{\sqrt{1 + \log i}} > t\right) dt \\ &\leq 2\sigma + \int_{2\sigma}^\infty \mathbb{P}\left(\max_{i \in \{1, \dots, N\}} \frac{|X_i|}{\sqrt{1 + \log i}} > t\right) dt \\ &\leq 2\sigma + \int_{2\sigma}^\infty 4 \exp\left(-\frac{t^2}{2\sigma^2}\right) dt \\ &\leq 3\sigma. \end{aligned}$$

Then we have

$$\mathbb{E} \max_{i \in \{1, \dots, N\}} |X_i| \leq \sqrt{1 + \log N} \left[\mathbb{E} \max_{i \in \{1, \dots, N\}} \frac{|X_i|}{\sqrt{1 + \log i}} \right] \leq \sqrt{1 + \log N} 3\sigma = \sqrt{9\sigma^2 \log N}.$$

□

Lemma E.7. *Suppose $\{X^{(k)}, k = 1, 2, \dots\}$ is a sequence of random vectors in \mathbb{R}^d and $X^{(k)} \rightsquigarrow X$ in \mathbb{R}^d . Then for d fixed elements v_1, \dots, v_d in a normed vector space B , we have $\sum_{i=1}^d X_i^{(k)} v_i \rightsquigarrow \sum_{i=1}^d X_i v_i$ in B .*

Proof. It suffices to verify that for any bounded continuous functions $f: B \rightarrow \mathbb{R}$, we have

$$\mathbb{E} f \left(\sum_{i=1}^d X_i^{(k)} v_i \right) \rightarrow \mathbb{E} f \left(\sum_{i=1}^d X_i v_i \right).$$

We may define $g(x) = f \left(\sum_{i=1}^d x_i v_i \right)$ as a function on \mathbb{R}^d . Since f is continuous, we have $\forall \epsilon > 0$, $\exists \delta > 0$ such that $\|u - v\| \leq \delta$ implies $|f(u) - f(v)| \leq \epsilon$. Let $V := \max_{i \in \{1, \dots, d\}} \|v_i\|$, we have $\left\| \sum_{i=1}^d x_i^{(1)} v_i - \sum_{i=1}^d x_i^{(2)} v_i \right\| \leq \delta$ as long as $\|x^{(1)} - x^{(2)}\|_1 \leq \delta/V$. Therefore, we may conclude g is a bounded continuous function on \mathbb{R}^d , and

$$\mathbb{E} f \left(\sum_{i=1}^d X_i^{(k)} v_i \right) = \mathbb{E} g \left(X^{(k)} \right) \rightarrow \mathbb{E} g \left(X \right) = \mathbb{E} f \left(\sum_{i=1}^d X_i v_i \right)$$

due to $X^{(k)} \rightsquigarrow X$. □

Lemma E.8. *Let $F^{-1}(p) := \inf \{t \mid F(t) \geq p\}$ be the quantile function of cumulative distribution function F . For any sequence of cumulative distribution functions, $F_n^{-1} \rightsquigarrow F^{-1}$ if and only if $F_n \rightsquigarrow F$. Here $F_n^{-1} \rightsquigarrow F^{-1}$ means that $F_n^{-1}(p) \rightsquigarrow F^{-1}(p)$ for every p where F^{-1} is continuous.*

Proof. See Lemma 21.2 in [40]. □

Lemma E.9. *Assume \mathcal{X} is a normed vector space and $g(x, z): \mathcal{X} \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a function such that for any fixed $z_0 \in \mathbb{R}^d$, $g(x, z_0)$ is continuous in x at x_0 . Suppose Z is a random vector taking values in \mathbb{R}^d , $F_x(t) := \mathbb{P}(g(x, Z) \leq t)$ is the distribution function of $g(x, Z)$ and $F_x^{-1}(p) := \inf \{t \mid F_x(t) \geq p\}$ is the quantile function of $g(x, Z)$. Then we have*

1. $F_x(t)$ is continuous in x at x_0 whenever t is a continuous point of $F_{x_0}(t)$;
2. $F_x^{-1}(p)$ is continuous in x at x_0 whenever p is a continuous point of $F_{x_0}^{-1}(p)$.

Proof. For any sequence $\{x_n\}$ such that $x_n \rightarrow x_0$, we have $g(x_n, Z) \xrightarrow{a.s.} g(x_0, Z)$, thus $F_{x_n} \rightsquigarrow F_{x_0}$, which further implies $F_x(t)$ is continuous at x_0 for every t where F_{x_0} is continuous. Then we may use Lemma E.8 to get $F_{x_n}^{-1} \rightsquigarrow F_{x_0}^{-1}$, which implies $F_{x_n}^{-1}(p) \rightarrow F_{x_0}^{-1}(p)$ whenever p is a continuous point of F^{-1} . Therefore, our conclusion follows. \square