

# A Finite Sample Analysis of Distributional TD Learning with Linear Function Approximation

Yang Peng<sup>\*</sup>      Kaicheng Jin<sup>†</sup>      Liangyu Zhang,<sup>‡</sup>      Zhihua Zhang<sup>§</sup>

May 14, 2025

## Abstract

In this paper, we study the finite-sample statistical rates of distributional temporal difference (TD) learning with linear function approximation. The aim of distributional TD learning is to estimate the return distribution of a discounted Markov decision process for a given policy  $\pi$ . Previous works on statistical analysis of distributional TD learning mainly focus on the tabular case. In contrast, we first consider the linear function approximation setting and derive sharp finite-sample rates. Our theoretical results demonstrate that the sample complexity of linear distributional TD learning matches that of classic linear TD learning. This implies that, with linear function approximation, learning the full distribution of the return from streaming data is no more difficult than learning its expectation (value function). To derive tight sample complexity bounds, we conduct a fine-grained analysis of the linear-categorical Bellman equation and employ the exponential stability arguments for products of random matrices. Our results provide new insights into the statistical efficiency of distributional reinforcement learning algorithms.

## 1 Introduction

Distributional policy evaluation [36, 2, 4], which aims to estimate the return distribution of a policy in an Markov decision process (MDP), is crucial for many uncertainty-aware or risk-sensitive tasks [30, 21]. Unlike the classic policy evaluation that only focuses on expected returns (value functions), distributional policy evaluation captures uncertainty and risk by considering the full distributional

---

<sup>\*</sup>School of Mathematical Sciences, Peking University; email: pengyang@pku.edu.cn.

<sup>†</sup>School of Mathematical Sciences, Peking University; email: kcjin@pku.edu.cn.

<sup>‡</sup>School of Statistics and Data Science, Shanghai University of Finance and Economics; email: zhangliangyu@sufe.edu.cn.

<sup>§</sup>School of Mathematical Sciences, Peking University; email: zhzhzhang@math.pku.edu.cn.

information. To solve a distributional policy evaluation problem, in the seminal work [2] proposed distributional temporal difference (TD) learning, which can be viewed as an extension of classic TD learning [58].

Although classic TD learning has been extensively studied [5, 61, 6, 14, 41, 25, 26, 10, 53, 54, 69], the theoretical understanding of distributional TD learning remains relatively underdeveloped. Recent works [48, 7, 73, 50, 51, 43] have analyzed distributional TD learning (or its model-based variants) in the tabular setting. Especially, Rowland et al. [51] and Peng et al. [43] demonstrated that in the tabular setting, learning the return distribution (in terms of the 1-Wasserstein distance<sup>1</sup>) is statistically as easy as learning its expectation. However, in practical scenarios, where the state space is extremely large or continuous, the function approximation [13, 12, 70, 17, 71, 39, 74, 31, 65, 57] becomes indispensable. This raises a new open question: *When function approximation is employed, does learning the return distribution remain as statistically efficient as learning its expectation?*

To answer this question, we consider the simplest form of function approximation, *i.e.*, linear function approximation, and investigate the finite-sample performance of linear distributional TD learning. In distributional TD learning, we need to represent the infinite-dimensional return distributions with some finite-dimensional parametrizations to make the algorithm tractable. Previous works [3, 32, 4] have proposed various linear distributional TD learning algorithms under different parameterizations, namely categorical and quantile parametrizations. In this paper, we consider the categorical parametrization and propose an improved version of the linear-categorical TD learning algorithm (**Linear-CTD**). We then analyze the non-asymptotic convergence rate of **Linear-CTD**. Our analysis reveals that, with the Polyak-Ruppert tail averaging [52, 45] and a proper constant step size, the sample complexity of **Linear-CTD** matches that of classic linear TD learning (**Linear-TD**) [26, 54]. Thus, this confirms that learning the return distribution is statistically no more difficult than learning its expectation when the linear function approximation is employed.

**Notation.** In the following parts of the paper,  $(x)_+ := \max\{x, 0\}$  for any  $x \in \mathbb{R}$ . “ $\lesssim$ ” (resp. “ $\gtrsim$ ”) means no larger (resp. smaller) than up to a multiplicative universal constant, and  $a \simeq b$  means  $a \lesssim b$  and  $a \gtrsim b$  hold simultaneously. The asymptotic notation  $f(\cdot) = \tilde{O}(g(\cdot))$  (resp.  $\tilde{\Omega}(g(\cdot))$ ) means that  $f(\cdot)$  is order-wise no larger (resp. smaller) than  $g(\cdot)$ , ignoring logarithmic factors of polynomials of  $(1 - \gamma)^{-1}, \lambda_{\min}^{-1}, \alpha^{-1}, \varepsilon^{-1}, \delta^{-1}, K, \|\psi^*\|_{\Sigma_\phi}, \|\theta^*\|_{I_K \otimes \Sigma_\phi}$ . We will explain the concrete meaning of

---

<sup>1</sup>Solving distributional policy evaluation  $\varepsilon$ -accurately in the 1-Wasserstein distance sense is harder than solving classic policy evaluation  $\varepsilon$ -accurately, as the absolute difference of value functions is always bounded by the 1-Wasserstein distance between return distributions.

the notation once we have encountered them for the first time.

We denote by  $\delta_x$  the Dirac measure at  $x \in \mathbb{R}$ ,  $\mathbb{1}$  the indicator function,  $\otimes$  the Kronecker product (see Appendix A),  $\mathbf{1}_K \in \mathbb{R}^K$  the all-ones vector,  $\mathbf{0}_K \in \mathbb{R}^K$  the all-zeros vector,  $\mathbf{I}_K \in \mathbb{R}^{K \times K}$  the identity matrix,  $\|\mathbf{u}\|$  the Euclidean norm of any vector  $\mathbf{u}$ ,  $\|\mathbf{B}\|$  the spectral norm of any matrix  $\mathbf{B}$ , and  $\|\mathbf{u}\|_B := \sqrt{\mathbf{u}^\top \mathbf{B} \mathbf{u}}$  when  $\mathbf{B}$  is positive semi-definite (PSD).  $\mathbf{B}_1 \preceq \mathbf{B}_2$  stands for  $\mathbf{B}_2 - \mathbf{B}_1$  is PSD for any symmetric matrices  $\mathbf{B}_1, \mathbf{B}_2$ . And  $\prod_{k=1}^t \mathbf{B}_k$  is defined as  $\mathbf{B}_t \mathbf{B}_{t-1} \cdots \mathbf{B}_1$  for any matrices  $\{\mathbf{B}_k\}_{k=1}^t$  with appropriate sizes. For any matrix  $\mathbf{B} = [\mathbf{b}(1), \dots, \mathbf{b}(n)] \in \mathbb{R}^{m \times n}$ , we define its vectorization as  $\text{vec}(\mathbf{B}) = (\mathbf{b}(1)^\top, \dots, \mathbf{b}(n)^\top)^\top \in \mathbb{R}^{mn}$ . Given a set  $A$ , we denote by  $\Delta(A)$  the set of all probability distributions over  $A$ . For simplicity, we abbreviate  $\Delta([0, (1-\gamma)^{-1}])$  as  $\mathcal{P}$ .

**Contributions.** Our contribution is two-fold: in algorithms and in theory. Algorithmically, we propose an improved version of the linear-categorical TD learning algorithm (**Linear-CTD**). Rather than using stochastic semi-gradient descent to update the parameter as in Bellemare et al. [3], Lyle et al. [32], Bellemare et al. [4], we directly formulate the linear-categorical projected Bellman equation into a linear system and apply a linear stochastic approximation to solve it. The resulting **Linear-CTD** can be viewed as a preconditioned version [9, 29] of the vanilla linear categorical TD learning algorithm proposed in Bellemare et al. [4, Section 9.6]. By introducing a preconditioner, our **Linear-CTD** achieves a finite-sample rate independent of the number of supports  $K$  in the categorical parameterization, which the vanilla version cannot attain. We provide both theoretical and experimental evidence to demonstrate this advantage of our **Linear-CTD**.

Theoretically, we establish the first non-asymptotic guarantees for distributional TD learning with the linear function approximation. Specifically, we show that in the generative model setting, with the Polyak-Ruppert tail averaging and a constant step size, we need

$$T = \tilde{O} \left( (\varepsilon^{-2} + \lambda_{\min}^{-1}) (1 - \gamma)^{-2} \lambda_{\min}^{-1} \left( K^{-1} (1 - \gamma)^{-2} \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi}^2 + 1 \right) \right)$$

online interactions with the environment to ensure **Linear-CTD** yields a  $\varepsilon$ -accurate estimator with high probability, when the error is measured by the  $\mu_\pi$ -weighted 1-Wasserstein distance. We also extend the result to the Markovian setting. Our sample complexity bounds match those of the classic **Linear-TD** with a constant step size [26, 54], confirming the same statistical tractability of distributional and classic value-based policy evaluations. To establish these theoretical results, we analyze the linear-categorical Bellman equation in detail and apply the exponential stability argument proposed in Samsonov et al. [54]. Our analysis of the linear-categorical Bellman equation lays the

foundation for subsequent algorithmic and theoretical advances in distributional reinforcement learning with function approximation.

**Organization.** The remainder of this paper is organized as follows. In Section 2, we recap **Linear**-TD and tabular categorical TD learning. In Section 3, we introduce the linear-categorical parametrization, and use the linear-categorical projected Bellman equation to derive **Linear**-CTD. In Section 4, we employ the exponential stability arguments to analyze the statistical efficiency of **Linear**-CTD. The proof is outlined in Section 5. In Section 6, we conclude our work. See Appendix B for more related work. In Appendix F, we empirically validate the convergence of **Linear**-CTD and compare it with prior algorithms through numerical experiments, confirming our theoretical findings. Details of the proof are given in the appendices.

## 2 Backgrounds

In this section, we recap the basics of policy evaluation and distributional policy evaluation tasks.

### 2.1 Policy Evaluation

A discounted MDP is defined by a 4-tuple  $M = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \gamma \rangle$ . We assume that the state space  $\mathcal{S}$  and the action space  $\mathcal{A}$  are both Polish spaces, namely complete separable metric spaces.  $\mathcal{P}(\cdot, \cdot \mid s, a)$  is the joint distribution of reward and next state condition on  $(s, a) \in \mathcal{S} \times \mathcal{A}$ . We assume that all rewards are bounded random variables in  $[0, 1]$ . And  $\gamma \in (0, 1)$  is the discount factor.

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:  $a_t \mid s_t \sim \pi(\cdot \mid s_t)$ ,  $(r_t, s_{t+1}) \mid (s_t, a_t) \sim \mathcal{P}(\cdot, \cdot \mid s_t, a_t)$ , for any  $t \in \mathbb{N}$ . We assume the Markov chain  $\{s_t\}_{t=0}^{\infty}$  has a unique stationary distribution  $\mu_{\pi} \in \Delta(\mathcal{S})$ . We define the return of the trajectory as  $G^{\pi}(s) := \sum_{t=0}^{\infty} \gamma^t r_t$ . The value function  $V^{\pi}(s)$  is the expectation of  $G^{\pi}(s)$ , and  $\mathbf{V}^{\pi} := (V^{\pi}(s))_{s \in \mathcal{S}} \in \mathbb{R}^{\mathcal{S}}$ . It is known that  $\mathbf{V}^{\pi}$  satisfies the Bellman equation:

$$V^{\pi}(s) = \mathbb{E}_{a \sim \pi(\cdot \mid s), (r, s') \sim \mathcal{P}(\cdot, \cdot \mid s, a)}[r + \gamma V^{\pi}(s')], \quad \forall s \in \mathcal{S}, \quad (1)$$

or in a compact form  $\mathbf{V}^{\pi} = \mathbf{T}^{\pi} \mathbf{V}^{\pi}$ , where  $\mathbf{T}^{\pi}: \mathbb{R}^{\mathcal{S}} \rightarrow \mathbb{R}^{\mathcal{S}}$  is called the Bellman operator. In the task of policy evaluation, we aim to find the unique solution  $\mathbf{V}^{\pi}$  of the equation for some given policy  $\pi$ .

**Tabular TD Learning.** The policy evaluation problem is reduced to solving the Bellman equation. However, in practical applications  $\mathbf{T}^\pi$  is usually unknown and the agent only has access to the streaming data  $\{(s_t, a_t, r_t)\}_{t=0}^\infty$ . In this circumstance, we can solve the Bellman equation through linear stochastic approximation (LSA). Specifically, in the  $t$ -th time-step the updating scheme is

$$V_t(s_t) \leftarrow V_{t-1}(s_t) - \alpha (V_{t-1}(s_t) - r_t - \gamma V_{t-1}(s_{t+1})), \quad V_t(s) \leftarrow V_{t-1}(s), \quad \forall s \neq s_t. \quad (2)$$

We expect  $\mathbf{V}_t$  to converge to  $\mathbf{V}^\pi$  as  $t$  tends to infinity. This algorithm is known as TD learning, however, it is computationally tractable only in the tabular setting.

**Linear Function Approximation and Linear TD Learning.** In this part, we introduce linear function approximation and briefly review the more practical **Linear-TD**. To be concrete, we assume there is a  $d$ -dimensional feature vector for each state  $s \in \mathcal{S}$ , which is given by the feature map  $\phi: \mathcal{S} \rightarrow \mathbb{R}^d$ . We consider the linear function approximation of value functions:

$$\mathcal{V}_\phi := \left\{ \mathbf{V}_\psi = (V_\psi(s))_{s \in \mathcal{S}} : V_\psi(s) = \phi(s)^\top \psi, \psi \in \mathbb{R}^d \right\} \subset \mathbb{R}^\mathcal{S}, \quad (3)$$

$\mu_\pi$ -weighted norm  $\|\mathbf{V}\|_{\mu_\pi} := (\mathbb{E}_{s \sim \mu_\pi} [V(s)^2])^{1/2}$ , and linear projection operator  $\Pi_\phi^\pi: \mathbb{R}^\mathcal{S} \rightarrow \mathcal{V}_\phi$ :

$$\Pi_\phi^\pi \mathbf{V} := \operatorname{argmin}_{\mathbf{V}_\psi \in \mathcal{V}_\phi} \|\mathbf{V} - \mathbf{V}_\psi\|_{\mu_\pi}, \quad \forall \mathbf{V} \in \mathbb{R}^\mathcal{S}.$$

One can check that the linear projected Bellman operator  $\Pi_\phi^\pi \mathbf{T}^\pi$  is a  $\gamma$ -contraction in the Polish space  $(\mathcal{V}_\phi, \|\cdot\|_{\mu_\pi})$ . Hence,  $\Pi_\phi^\pi \mathbf{T}^\pi$  admits a unique fixed point  $\mathbf{V}_{\psi^\star}$ , which satisfies  $\|\mathbf{V}^\pi - \mathbf{V}_{\psi^\star}\|_{\mu_\pi} \leq (1-\gamma^2)^{-1/2} \|\mathbf{V}^\pi - \Pi_\phi^\pi \mathbf{V}^\pi\|_{\mu_\pi}$  [4, Theorem 9.8]. In Appendix C.1, we show that  $\psi^\star \in \mathbb{R}^d$  is the unique solution to the linear system for  $\psi \in \mathbb{R}^d$ :

$$(\Sigma_\phi - \gamma \mathbb{E}_{s,s'} [\phi(s)\phi(s')^\top]) \psi = \mathbb{E}_{s,r} [\phi(s)r], \quad \Sigma_\phi := \mathbb{E}_{s \sim \mu_\pi} [\phi(s)\phi(s)^\top]. \quad (4)$$

In the subscript of the expectation, we abbreviate  $s \sim \mu_\pi(\cdot), a \sim \pi(\cdot|s), (r, s') \sim \mathcal{P}(\cdot, \cdot | s, a)$  as  $s, a, r, s'$ . For brevity, we will use such abbreviations in this paper when there is no ambiguity. We can use LSA to solve the linear projected Bellman equation (Eqn. 4). As a result, at the  $t$ -th time-step, the updating scheme of **Linear-TD** is

$$\text{Linear-TD:} \quad \psi_t \leftarrow \psi_{t-1} - \alpha \phi(s_t) \left[ (\phi(s_t) - \gamma \phi(s_{t+1}))^\top \psi_{t-1} - r_t \right]. \quad (5)$$

## 2.2 Distributional Policy Evaluation

In certain applications, we are not only interested in finding the expectation of random return  $G^\pi(s)$  but also want to find the whole distribution of  $G^\pi(s)$ . This task is called distributional policy evaluation. We use  $\eta^\pi(s) \in \mathcal{P}$  to denote the distribution of  $G^\pi(s)$  and let  $\boldsymbol{\eta}^\pi := (\eta^\pi(s))_{s \in \mathcal{S}} \in \mathcal{P}^{\mathcal{S}}$ . Then  $\boldsymbol{\eta}^\pi$  satisfies the distributional Bellman equation:

$$\eta^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s), (r, s') \sim \mathcal{P}(\cdot, \cdot|s, a)}[(b_{r, \gamma})_\# \eta^\pi(s')], \quad \forall s \in \mathcal{S}, \quad (6)$$

where the RHS is the distribution of  $r_0 + \gamma G^\pi(s_1)$  conditioned on  $s_0 = s$ . Here  $b_{r, \gamma}(x) := r + \gamma x$  for any  $x \in \mathbb{R}$ , and  $f_\# \nu \in \mathcal{P}$  is defined as  $f_\# \nu(A) := \nu(\{x: f(x) \in A\})$  for any function  $f: \mathbb{R} \rightarrow \mathbb{R}$ , probability measure  $\nu \in \mathcal{P}$  and Borel set  $A \subset \mathbb{R}$ . The distributional Bellman equation can also be written as  $\boldsymbol{\eta}^\pi = \boldsymbol{\mathcal{T}}^\pi \boldsymbol{\eta}^\pi$ . The operator  $\boldsymbol{\mathcal{T}}^\pi: \mathcal{P}^{\mathcal{S}} \rightarrow \mathcal{P}^{\mathcal{S}}$  is called the distributional Bellman operator. In this task, our goal is to find  $\boldsymbol{\eta}^\pi$  for some given policy  $\pi$ .

**Tabular Distributional TD Learning.** In analogy to tabular TD learning (Eqn. (2)), in the tabular setting, we can solve the distributional Bellman equation by LSA and derive the distributional TD learning rule given the streaming data  $\{(s_t, a_t, r_t)\}_{t=0}^\infty$ :

$$\eta_t(s_t) \leftarrow \eta_{t-1}(s_t) - \alpha[\eta_{t-1}(s_t) - (b_{r_t, \gamma})_\# \eta_{t-1}(s_{t+1})], \quad \eta_t(s) \leftarrow \eta_{t-1}(s), \quad \forall s \neq s_t.$$

We comment the algorithm above is not computationally feasible as we need to manipulate infinite-dimensional objects (return distributions) at each iteration.

**Categorical Parametrization and Tabular Categorical TD Learning.** In order to deal with return distributions in a computationally tractable manner, we consider the categorical parametrization as in Bellemare et al. [2], Rowland et al. [48, 51], Peng et al. [43]. To be compatible with linear function approximation introduced in the next section, which cannot guarantee non-negative outputs, we will work with  $\mathcal{P}^{\text{sign}}$ , the signed measure space with total mass 1 as in Bellemare et al. [3], Lyle et al. [32], Bellemare et al. [4] instead of standard probability space  $\mathcal{P} \subset \mathcal{P}^{\text{sign}}$ :

$$\mathcal{P}^{\text{sign}} := \{\nu: \nu(\mathbb{R}) = 1, \text{supp}(\nu) \subseteq [0, (1 - \gamma)^{-1}]\}.$$

For any  $\nu \in \mathcal{P}^{\text{sign}}$ , we define its cumulative distribution function (CDF) as  $F_\nu(x) := \nu([0, x])$ . We can naturally define the  $L^2$  and  $L^1$  distances between CDFs as the Cramér distance  $\ell_2$  and 1-Wasserstein distance  $W_1$  in  $\mathcal{P}^{\text{sign}}$ , respectively. The distributional Bellman operator (see Eqn. (6)) can also be extended to the product space  $(\mathcal{P}^{\text{sign}})^{\mathcal{S}}$  without modifying its definition.

The space of all categorical parametrized signed measures with total mass 1 is defined as

$$\mathcal{P}_K^{\text{sign}} := \left\{ \nu_{\mathbf{p}} = \sum_{k=0}^K p_k \delta_{x_k} : \mathbf{p} = (p_0, \dots, p_{K-1})^\top \in \mathbb{R}^K, p_K = 1 - \sum_{k=0}^{K-1} p_k \right\}, \quad (7)$$

which is an affine subspace of  $\mathcal{P}^{\text{sign}}$ . Here  $\{x_k = k\iota_K\}_{k=0}^K$  are  $K+1$  equally-spaced points of the support,  $\iota_K = [K(1-\gamma)]^{-1}$  is the gap between adjacent points, and  $p_k$  is the ‘probability’ (may be negative) that  $\nu$  assigns to  $x_k$ . We define the categorical projection operator  $\Pi_K : \mathcal{P}^{\text{sign}} \rightarrow \mathcal{P}_K^{\text{sign}}$  as

$$\Pi_K \nu := \operatorname{argmin}_{\nu_{\mathbf{p}} \in \mathcal{P}_K^{\text{sign}}} \ell_2(\nu, \nu_{\mathbf{p}}), \quad \forall \nu \in \mathcal{P}^{\text{sign}}.$$

Following Bellemare et al. [4, Proposition 5.14], one can show that  $\Pi_K \nu \in \mathcal{P}_K^{\text{sign}}$  is uniquely represented with a vector  $\mathbf{p}_\nu = (p_k(\nu))_{k=0}^{K-1} \in \mathbb{R}^K$ , where

$$p_k(\nu) = \int_{[0, (1-\gamma)^{-1}]} (1 - |(x - x_k)/\iota_K|)_+ \nu(dx). \quad (8)$$

We lift  $\Pi_K$  to the product space by defining  $[\Pi_K \boldsymbol{\eta}](s) := \Pi_K \eta(s)$ . One can check that the categorical Bellman operator  $\Pi_K \mathcal{T}^\pi$  is a  $\sqrt{\gamma}$ -contraction in the Polish space  $((\mathcal{P}_K^{\text{sign}})^{\mathcal{S}}, \ell_{2, \mu_\pi})$ , where  $\ell_{2, \mu_\pi}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2) := (\mathbb{E}_{s \sim \mu_\pi} [\ell_2^2(\eta_1(s), \eta_2(s))])^{1/2}$  is the  $\mu_\pi$ -weighted Cramér distance between  $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2 \in (\mathcal{P}^{\text{sign}})^{\mathcal{S}}$ . Similarly,  $W_{1, \mu_\pi}(\boldsymbol{\eta}_1, \boldsymbol{\eta}_2) := (\mathbb{E}_{s \sim \mu_\pi} [W_1^2(\eta_1(s), \eta_2(s))])^{1/2}$ . Hence, the categorical projected Bellman equation  $\boldsymbol{\eta} = \Pi_K \mathcal{T}^\pi \boldsymbol{\eta}$  admits a unique solution  $\boldsymbol{\eta}^{\pi, K}$ , which satisfies  $W_{1, \mu_\pi}(\boldsymbol{\eta}^\pi, \boldsymbol{\eta}^{\pi, K}) \leq (1-\gamma)^{-1} \ell_{2, \mu_\pi}(\boldsymbol{\eta}^\pi, \Pi_K \boldsymbol{\eta}^\pi)$  [48, Proposition 3]. Applying LSA to solving the equation yields tabular categorical TD learning, and the iteration rule is given by

$$\eta_t(s_t) \leftarrow \eta_{t-1}(s_t) - \alpha [\eta_{t-1}(s_t) - \Pi_K(b_{r_t, \gamma})_\# \eta_{t-1}(s_{t+1})], \quad \eta_t(s) \leftarrow \eta_{t-1}(s), \quad \forall s \neq s_t. \quad (9)$$

### 3 Linear-Categorical TD Learning

In this section, we propose our **Linear-CTD** algorithm (Eqn. (13)) by combining the linear function approximation (Eqn. (3)) with the categorical parametrization (Eqn. (7)). We first introduce the

space of linear-categorical parametrized signed measures with total mass 1:

$$\mathcal{P}_{\phi,K}^{\text{sign}} := \left\{ \boldsymbol{\eta}_{\boldsymbol{\theta}} = (\eta_{\boldsymbol{\theta}}(s))_{s \in \mathcal{S}} : \eta_{\boldsymbol{\theta}}(s) = \sum_{k=0}^K p_k(s; \boldsymbol{\theta}) \delta_{x_k}, \boldsymbol{\theta} = (\boldsymbol{\theta}(0)^\top, \dots, \boldsymbol{\theta}(K-1)^\top)^\top \in \mathbb{R}^{dK} \right\},$$

which is an affine subspace of  $(\mathcal{P}_K^{\text{sign}})^{\mathcal{S}}$ . Here  $p_k(s; \boldsymbol{\theta}) = F_k(s; \boldsymbol{\theta}) - F_{k-1}(s; \boldsymbol{\theta})$ , and

$$F_k(s; \boldsymbol{\theta}) = \boldsymbol{\phi}(s)^\top \boldsymbol{\theta}(k) + (k+1)/(K+1) \quad \text{for } k \in \{0, 1, \dots, K-1\} \quad (10)$$

is CDF of  $\eta_{\boldsymbol{\theta}}(s)$  at  $x_k$  ( $F_{-1}(s; \cdot) \equiv 0, F_K(s; \cdot) \equiv 1$ )<sup>2</sup>. In many cases, especially when formulating and implementing algorithms, it is much more convenient and efficient to work with the matrix version of the parameter  $\boldsymbol{\Theta} := (\boldsymbol{\theta}(0), \dots, \boldsymbol{\theta}(K-1)) \in \mathbb{R}^{d \times K}$  rather than with  $\boldsymbol{\theta} = \text{vec}(\boldsymbol{\Theta})$ . We define the linear-categorical projection operator  $\Pi_{\phi,K}^\pi : (\mathcal{P}_K^{\text{sign}})^{\mathcal{S}} \rightarrow \mathcal{P}_{\phi,K}^{\text{sign}}$  as follows:

$$\Pi_{\phi,K}^\pi \boldsymbol{\eta} := \underset{\boldsymbol{\eta}_{\boldsymbol{\theta}} \in \mathcal{P}_{\phi,K}^{\text{sign}}}{\text{argmin}} \ell_{2,\mu_\pi}(\boldsymbol{\eta}, \boldsymbol{\eta}_{\boldsymbol{\theta}}), \quad \forall \boldsymbol{\eta} \in (\mathcal{P}_K^{\text{sign}})^{\mathcal{S}}.$$

$\Pi_{\phi,K}^\pi$  is in fact an orthogonal projection (see Proposition D.2), and thus is non-expansive. The following proposition characterizes  $\Pi_{\phi,K}^\pi$ , whose proof can be found in Appendix D.2.

**Proposition 3.1.** *For any  $\boldsymbol{\eta} \in (\mathcal{P}_K^{\text{sign}})^{\mathcal{S}}$ ,  $\Pi_{\phi,K}^\pi \boldsymbol{\eta}$  is uniquely given by  $\boldsymbol{\eta}_{\tilde{\boldsymbol{\theta}}}$ , where  $\tilde{\boldsymbol{\theta}} = \text{vec}(\tilde{\boldsymbol{\Theta}})$ ,*

$$\tilde{\boldsymbol{\Theta}} = \boldsymbol{\Sigma}_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} [\boldsymbol{\phi}(s)(\mathbf{p}_\eta(s) - (K+1)^{-1} \mathbf{1}_K)^\top \mathbf{C}^\top], \quad \mathbf{C} = [\mathbf{1}\{i \geq j\}]_{i,j \in [K]} \in \mathbb{R}^{K \times K}. \quad (11)$$

Here  $\mathbf{p}_\eta(s) := \mathbf{p}_{\eta(s)} = (p_k(\eta(s)))_{k=0}^{K-1}$  is the vector that identifies  $\Pi_K \eta(s)$  defined in Eqn. (8).

Since  $\Pi_{\phi,K}^\pi \mathcal{T}^\pi$  is a  $\sqrt{\gamma}$ -contraction in  $(\mathcal{P}_{\phi,K}^{\text{sign}}, \ell_{2,\mu_\pi})$  ( $\mathcal{T}^\pi$  is  $\sqrt{\gamma}$ -contraction [4, Lemma 9.14]), in the following theorem, we can generalize the linear projected Bellman equation (Eqn. (4)) to the distributional setting. The proof can be found in Appendix D.3.

**Theorem 3.1.** *The linear-categorical projected Bellman equation  $\boldsymbol{\eta}_{\boldsymbol{\theta}} = \Pi_{\phi,K}^\pi \mathcal{T}^\pi \boldsymbol{\eta}_{\boldsymbol{\theta}}$  admits a unique solution  $\boldsymbol{\eta}_{\boldsymbol{\theta}^*}$ , where the matrix parameter  $\boldsymbol{\Theta}^*$  is the unique solution to the linear system for  $\boldsymbol{\Theta} \in \mathbb{R}^{d \times K}$*

$$\boldsymbol{\Sigma}_\phi \boldsymbol{\Theta} - \mathbb{E}_{s,s',r} \left[ \boldsymbol{\phi}(s) \boldsymbol{\phi}(s')^\top \boldsymbol{\Theta} (\mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1})^\top \right] = \frac{1}{K+1} \mathbb{E}_{s,r} \left[ \boldsymbol{\phi}(s) \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right)^\top \mathbf{C}^\top \right], \quad (12)$$

<sup>2</sup>The  $(k+1)/(K+1)$  term in CDF (Eqn. (10)) corresponds to a discrete uniform distribution  $\nu$  for normalization, i.e., making sure  $F_K(s; \cdot) \equiv 1$ . If we include an intercept term in the feature or we consider the tabular setting, we can replace  $\nu$  with any distribution without affecting the definition of  $\mathcal{P}_{\phi,K}^{\text{sign}}$  or any other things.



where for any  $r \in [0, 1]$  and  $j, k \in \{0, 1, \dots, K\}$ ,

$$g_{j,k}(r) := (1 - |(r + \gamma x_j - x_k)/\iota_K|)_+, \quad \mathbf{g}_j(r) := (g_{j,k}(r))_{k=0}^{K-1} \in \mathbb{R}^K,$$

$$\mathbf{G}(r) := [\mathbf{g}_0(r), \dots, \mathbf{g}_{K-1}(r)] \in \mathbb{R}^{K \times K}, \quad \tilde{\mathbf{G}}(r) := \mathbf{G}(r) - \mathbf{1}_K^\top \otimes \mathbf{g}_K(r) \in \mathbb{R}^{K \times K}.$$

In analogy to the approximation bounds of  $\|\mathbf{V}^\pi - \mathbf{V}_{\psi^\star}\|_{\mu_\pi}$  and  $W_{1,\mu_\pi}(\boldsymbol{\eta}^\pi, \boldsymbol{\eta}^{\pi,K})$ , the following lemma answers how close  $\boldsymbol{\eta}_{\theta^\star}$  is to  $\boldsymbol{\eta}^\pi$ , whose proof can be found in Appendix D.5.

**Proposition 3.2** (Approximation Error of  $\boldsymbol{\eta}_{\theta^\star}$ ). *It holds that*

$$W_{1,\mu_\pi}^2(\boldsymbol{\eta}^\pi, \boldsymbol{\eta}_{\theta^\star}) \leq K^{-1}(1 - \gamma)^{-3} + (1 - \gamma)^{-2} \ell_{2,\mu_\pi}^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_{\phi,K}^\pi \boldsymbol{\eta}^\pi),$$

where the first error term  $K^{-1}(1 - \gamma)^{-3}$  is due to the categorical parametrization, and the second error term  $(1 - \gamma)^{-2} \ell_{2,\mu_\pi}^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_{\phi,K}^\pi \boldsymbol{\eta}^\pi)$  is due to the additional linear function approximation.

As before, we solve Eqn. (12) by LSA and get **Linear-CTD** given the streaming data  $\{(s_t, a_t, r_t)\}_{t=0}^\infty$ :

$$\begin{aligned} \text{Linear-CTD:} \quad \boldsymbol{\Theta}_t \leftarrow & \boldsymbol{\Theta}_{t-1} - \alpha \phi(s_t) \left[ \phi(s_t)^\top \boldsymbol{\Theta}_{t-1} - \phi(s_{t+1})^\top \boldsymbol{\Theta}_{t-1} (\mathbf{C} \tilde{\mathbf{G}}(r_t) \mathbf{C}^{-1})^\top \right. \\ & \left. - (K+1)^{-1} (\sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K)^\top \mathbf{C}^\top \right], \end{aligned} \quad (13)$$

for any  $t \geq 1$ , where  $\alpha$  is the constant step size. In Appendix F, we empirically validate the convergence of **Linear-CTD** through numerical experiments. It is easy to verify that, in the special case of the tabular setting  $(\phi(s) = (\mathbb{1}\{s=\tilde{s}\})_{\tilde{s} \in \mathcal{S}})$ , **Linear-CTD** is equivalent to tabular categorical TD learning (Eqn. (9)). In this paper, we consider the Polyak-Ruppert tail averaging  $\bar{\boldsymbol{\theta}}_T := (T/2+1)^{-1} \sum_{t=T/2}^T \boldsymbol{\theta}_t$  (we use an even number  $T$ ) as in the analysis of **Linear-TD** in Samsonov et al. [54]. Standard theory of LSA [37] says under some conditions, if we take an appropriate step size  $\alpha$ ,  $\bar{\boldsymbol{\theta}}_T$  will converge to the solution  $\boldsymbol{\theta}^\star$  with rate  $T^{-1/2}$  as  $T \rightarrow \infty$ .

**Remark 1** (Comparison with Existing Linear Distributional TD Learning Algorithms). *Our **Linear-CTD** can be regarded as a preconditioned version of vanilla stochastic semi-gradient descent (SSGD) with the probability mass function (PMF) representation [4, Section 9.6]. See Appendix D.6 for the PMF representation, and Appendix D.7 for a self-contained derivation of SSGD with PMF representations. The preconditioning technique is a commonly used methodology to accelerate solving optimization problems by reducing the condition number. We precondition the vanilla algorithm by removing the matrix  $\mathbf{C}^\top \mathbf{C}$  (see Eqn. (26)), whose condition number scales with  $K^2$  (Lemma G.2).*

By introducing the preconditioner  $(\mathbf{C}^\top \mathbf{C})^{-1}$ , our **Linear-CTD** (Eqn. (13)) can achieve a convergence rate independent of  $K$ , which the vanilla form cannot achieve. See Remark 5 and Appendix F for theoretical and experimental evidence respectively.

We comment that **Linear-CTD** (Eqn. (13)) is equivalent to SSGD with CDF representation, which was also considered in Lyle et al. [32]. The difference is that our **Linear-CTD** normalizes the distribution so that the total mass of return distributions always be 1, while the algorithm in Lyle et al. [32] does not. See Appendix D.7 for a self-contained derivation of SSGD with CDF representations.

Bellemare et al. [4, Section 9.5] also proposed a softmax-based linear-categorical algorithm, which is closer to the practical C51 algorithm [2]. However, the nonlinearity due to softmax makes it difficult for analysis. We will leave the analysis for it as future work.

**Remark 2** (**Linear-CTD** is mean-preserving). A key property of **Linear-CTD** is mean preservation. That is, if we use identical initializations ( $\mathbb{E}_{G \sim \eta_{\theta_0}}[G] = \mathbf{V}_{\psi_0}$ ) and an identical data stream to update in both **Linear-CTD** and **Linear-TD**, it follows that  $\mathbb{E}_{G \sim \eta_{\theta_t}}[G] = \mathbf{V}_{\psi_t}$  for all  $t$ . However, the mean-preserving property does not hold for the SSGD with the PMF representation. In this sense, our **Linear-CTD** is indeed the generalization of **Linear-TD**. See Appendix D.8 for details.

## 4 Non-Asymptotic Statistical Analysis

In our task, the quality of estimator  $\eta_{\bar{\theta}_T}$  is measured by  $\mu_\pi$ -weighted 1-Wasserstein error  $W_{1,\mu_\pi}(\eta_{\bar{\theta}_T}, \eta^\pi)$ . By triangle inequality, the error can be decomposed into the approximation error and the estimation error:  $W_{1,\mu_\pi}(\eta^\pi, \eta_{\bar{\theta}_T}) \leq W_{1,\mu_\pi}(\eta^\pi, \eta_{\theta^\star}) + W_{1,\mu_\pi}(\eta_{\theta^\star}, \eta_{\bar{\theta}_T})$ . Proposition 3.2 already provided an upper bound for the approximation error  $W_{1,\mu_\pi}(\eta^\pi, \eta_{\theta^\star})$ , so it suffices to control the estimation error  $W_{1,\mu_\pi}(\eta_{\theta^\star}, \eta_{\bar{\theta}_T})$ , denoted  $\mathcal{L}(\bar{\theta}_T)$ .

In the following theorem, we give non-asymptotic convergence rates of  $\mathcal{L}(\bar{\theta}_T)$ . We start from the generative model setting, *i.e.*, in the  $t$ -th iteration, we collect samples  $s_t \sim \mu_\pi(\cdot)$ ,  $a_t \sim \pi(\cdot|s_t)$ ,  $(r_t, s'_t) \sim \mathcal{P}(\cdot, \cdot|s_t, a_t)$  from the generative model, and we replace  $s_{t+1}$  with  $s'_t$  in Eqn. (13). We give  $L^p$  and high-probability convergence results in this setting. These results can be extended to the Markovian setting, *i.e.*, using the streaming data  $\{(s_t, a_t, r_t)\}_{t=0}^\infty$ .

### 4.1 $L^2$ Convergence

We first provide non-asymptotic convergence rates of  $\mathbb{E}^{1/2}[(\mathcal{L}(\bar{\theta}_T))^2]$ , which do not grow with the number of supports  $K$ . The  $L^p$  ( $p > 2$ ) convergence results can be found in Theorem E.1.

**Theorem 4.1** ( $L^2$  Convergence). *For any  $K \geq (1 - \gamma)^{-1}$  and  $\alpha \in (0, (1 - \sqrt{\gamma})/76)$ , it holds that*

$$\begin{aligned} \mathbb{E}^{1/2}[(\mathcal{L}(\bar{\theta}_T))^2] &\lesssim \frac{\|\boldsymbol{\theta}^*\|_{V_1} + 1}{\sqrt{T}(1 - \gamma)\sqrt{\lambda_{\min}}} \left(1 + \sqrt{\frac{\alpha}{(1 - \gamma)\lambda_{\min}}}\right) + \frac{\|\boldsymbol{\theta}^*\|_{V_1} + 1}{T\sqrt{\alpha}(1 - \gamma)^{\frac{3}{2}}\lambda_{\min}} \\ &\quad + \frac{(1 - \frac{1}{2}\alpha(1 - \sqrt{\gamma})\lambda_{\min})^{T/2}}{T\sqrt{\alpha}(1 - \gamma)\sqrt{\lambda_{\min}}} \left(\frac{1}{\sqrt{\alpha}} + \frac{1}{\sqrt{(1 - \gamma)\lambda_{\min}}}\right) \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_{V_2}, \end{aligned}$$

where  $\|\boldsymbol{\theta}^*\|_{V_1} := \frac{1}{\sqrt{K}(1 - \gamma)} \|\boldsymbol{\theta}^*\|_{I_K \otimes \Sigma_\phi}$  and  $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_{V_2} := \frac{1}{\sqrt{K}(1 - \gamma)} \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|$ .

In the upper bound, the first term of order  $T^{-1/2}$  is dominant. The second term of order  $T^{-1}$  has a worse dependence on  $\alpha^{-1}$ , leading to a sample size barrier (Eqn. (15)). The third term, corresponding to the initialization error, decays at a geometric rate and can be thus ignored. To prove Theorem 4.1, we conduct a fine-grained analysis of the linear-categorical Bellman equation and apply the exponential stability argument [54]. We outline the proof in Section 5. In Remark 3, we compare our Theorem 4.1 with the  $L^2$  convergence rate of classic **Linear-TD** and conclude that learning the distribution of the return is as easy as learning its expectation (value function) with linear function approximation. Rowland et al. [51], Peng et al. [43] discovered this phenomenon in the tabular setting, and we extended it to the function approximation setting.

**Remark 3** (Comparison with Convergence Rate of **Linear-TD**). *The only difference between our Theorem 4.1 and the tight  $L^2$  convergence rate of classic **Linear-TD** (see Appendix C.2) lies in replacing  $\|\boldsymbol{\psi}^*\|_{\Sigma_\phi}$  (resp.  $\|\boldsymbol{\psi}_0 - \boldsymbol{\psi}^*\|$ ) in **Linear-TD** with  $\|\boldsymbol{\theta}^*\|_{V_1}$  (resp.  $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_{V_2}$ ). We claim that the two pairs should be of the same order, respectively. Note that  $\|\boldsymbol{\theta}^*\|_{V_1}$  and  $\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_{V_2}$  are of order  $\mathcal{O}((1 - \gamma)^{-1})$  if  $\eta_{\boldsymbol{\theta}^*}(s)$  and  $\eta_{\boldsymbol{\theta}_0}(s)$  are valid probability distributions for all  $s \in \mathcal{S}$ . This is because in this case,  $F_k(s; \boldsymbol{\theta}) = \boldsymbol{\phi}(s)^\top \boldsymbol{\theta}(k) + (k + 1)/(K + 1) \in [0, 1]$  for  $\boldsymbol{\theta} \in \{\boldsymbol{\theta}^*, \boldsymbol{\theta}_0\}$ . While in **Linear-TD**,  $\|\boldsymbol{\psi}^*\|_{\Sigma_\phi}$  and  $\|\boldsymbol{\psi}_0 - \boldsymbol{\psi}^*\|$  are also of order  $\mathcal{O}((1 - \gamma)^{-1})$  if  $V_\psi(s) = \boldsymbol{\phi}(s)^\top \boldsymbol{\psi} \in [0, (1 - \gamma)^{-1}]$  for all  $s \in \mathcal{S}$  and  $\boldsymbol{\psi} \in \{\boldsymbol{\psi}^*, \boldsymbol{\psi}_0\}$ . It is thus reasonable to consider the two pairs with the same order, respectively. Similar arguments also hold in other convergence results presented in this paper. Therefore, in this sense, our results match those of **Linear-TD**.*

One can translate Theorem 4.1 into a sample complexity bound.

**Corollary 4.1.** *Under the same conditions as in Theorem 4.1, for any  $\varepsilon > 0$ , suppose*

$$T \gtrsim \frac{\|\boldsymbol{\theta}^\star\|_{V_1}^2 + 1}{\varepsilon^2(1-\gamma)^2\lambda_{\min}} \left(1 + \frac{\alpha}{(1-\gamma)\lambda_{\min}}\right) + \frac{\|\boldsymbol{\theta}^\star\|_{V_1} + 1}{\varepsilon\sqrt{\alpha}(1-\gamma)^{\frac{3}{2}}\lambda_{\min}} \\ + \frac{1}{\alpha(1-\gamma)\lambda_{\min}} \left( \log \frac{\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\|_{V_2}}{\varepsilon} + \log \left( \frac{1}{T\sqrt{\alpha}(1-\gamma)\sqrt{\lambda_{\min}}} \left( \frac{1}{\sqrt{\alpha}} + \frac{1}{\sqrt{(1-\gamma)\lambda_{\min}}} \right) \right) \right).$$

*Then it holds that  $\mathbb{E}^{1/2}[(\mathcal{L}(\bar{\boldsymbol{\theta}}_T))^2] \leq \varepsilon$ .*

**Instance-Independent Step Size.** If we take the largest possible instance-independent step size, *i.e.*,  $\alpha \simeq (1-\gamma)$ , and consider  $\varepsilon \in (0, 1)$ , we obtain the sample complexity bound

$$T = \tilde{\mathcal{O}} \left( \varepsilon^{-2}(1-\gamma)^{-2}\lambda_{\min}^{-2} \left( \|\boldsymbol{\theta}^\star\|_{V_1}^2 + 1 \right) \right). \quad (14)$$

**Optimal Instance-Dependent Step Size.** If we take the optimal instance-dependent step size  $\alpha \simeq (1-\gamma)\lambda_{\min}$  which involves the unknown  $\lambda_{\min}$ , we obtain a better sample complexity bound

$$T = \tilde{\mathcal{O}} \left( (\varepsilon^{-2} + \lambda_{\min}^{-1}) (1-\gamma)^{-2}\lambda_{\min}^{-1} \left( \|\boldsymbol{\theta}^\star\|_{V_1}^2 + 1 \right) \right). \quad (15)$$

There is a sample size barrier in the bound, that is, the dependence on  $\lambda_{\min}$  is the optimal  $\lambda_{\min}^{-1}$  only when  $\varepsilon = \tilde{\mathcal{O}}(\sqrt{\lambda_{\min}})$ , or equivalently, we require a large enough (independent of  $\varepsilon$ ) update steps  $T$ .

These results match the recent results for classic **Linear-TD** with a constant step size [26, 54]. It is possible to break the sample size barrier in Eqn. (15) as in **Linear-TD** by applying variance-reduction techniques [27]. We leave it for future work.

## 4.2 Convergence with High Probability and Markovian Samples

Applying the  $L^p$  convergence result (Theorem E.1) with  $p = 2 \log(1/\delta)$  and Markov's inequality, we immediately obtain the high-probability convergence result.

**Theorem 4.2** (High-Probability Convergence). *For any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , suppose  $K \geq (1-\gamma)^{-1}$ ,  $\alpha \in (0, (1-\sqrt{\gamma})/[38 \log(T/\delta^2)])$ , and*

$$T = \tilde{\mathcal{O}} \left( \frac{\|\boldsymbol{\theta}^\star\|_{V_1}^2 + 1}{\varepsilon^2(1-\gamma)^2\lambda_{\min}} \left( 1 + \frac{\alpha \log \frac{1}{\delta}}{(1-\gamma)\lambda_{\min}} \right) \log \frac{1}{\delta} + \frac{\|\boldsymbol{\theta}^\star\|_{V_1} + 1}{\varepsilon\sqrt{\alpha}(1-\gamma)^{\frac{3}{2}}\lambda_{\min}} \log \frac{1}{\delta} + \frac{\log \frac{\|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\|_{V_2}}{\varepsilon}}{\alpha(1-\gamma)\lambda_{\min}} \right).$$

*Then with probability at least  $1-\delta$ , it holds that  $\mathcal{L}(\bar{\boldsymbol{\theta}}_T) \leq \varepsilon$ . Here, the  $\tilde{\mathcal{O}}(\cdot)$  does not hide polynomials*

of  $\log(1/\delta)$  (but hides logarithm terms of  $\log(1/\delta)$ ).

Again, we will obtain concrete sample complexity bounds as in Eqn. (14) or Eqn. (15) if we use different step sizes. Compared with the theoretical results for classic **Linear**-TD, our results match Samsonov et al. [54, Theorem 4]. [54] also considered the constant step size, but obtained a worse dependence on  $\log(1/\delta)$  than Wu et al. [69, Theorem 4] which uses the polynomial-decaying step size  $\alpha_t = \alpha_0 t^{-\beta}$  with  $\beta \in (1/2, 1)$  instead.

**Remark 4** (Markovian Setting). *Using the same argument as in the proof of Samsonov et al. [54, Theorem 6], one can immediately derive a high-probability sample complexity bound in the Markovian setting. Compared with the bound in the generative model setting (Theorem 4.2), the bound in the Markovian setting will have an additional dependency on  $t_{\text{mix}} \log(T/\delta)$ , where  $t_{\text{mix}}$  is the mixing time of the Markov chain  $\{s_t\}_{t=0}^{\infty}$  in  $\mathcal{S}$ . We omit this result for brevity.*

## 5 Proof Outlines

In this section, we outline the proofs of our main results (Theorem 4.1). We first state the theoretical properties of the linear-categorical Bellman equation and the exponential stability of **Linear**-CTD. Finally, we highlight some key steps in proving these results.

### 5.1 Vectorization of Linear-CTD

In our analysis, it will be more convenient to work with the vectorization version of the updating scheme of **Linear**-CTD (Eqn. (13)):

$$\begin{aligned} \boldsymbol{\theta}_t &\leftarrow \boldsymbol{\theta}_{t-1} - \alpha (\mathbf{A}_t \boldsymbol{\theta}_{t-1} - \mathbf{b}_t), \quad \mathbf{A}_t = [\mathbf{I}_K \otimes (\boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s_t)^\top)] - [(C \tilde{\mathbf{G}}(r_t) C^{-1}) \otimes (\boldsymbol{\phi}(s_t) \boldsymbol{\phi}(s'_t)^\top)], \\ \mathbf{b}_t &= (K+1)^{-1} \left[ C \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right) \right] \otimes \boldsymbol{\phi}(s_t). \end{aligned} \quad (16)$$

We denote by  $\bar{\mathbf{A}}$  and  $\bar{\mathbf{b}}$  the expectations of  $\mathbf{A}_t$  and  $\mathbf{b}_t$ , respectively. Using exponential stability arguments, we can derive an upper bound for  $\|\bar{\mathbf{A}}(\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*)\|$ . The following lemma further translates it to an upper bound for  $\mathcal{L}(\bar{\boldsymbol{\theta}}_T) = W_{1, \mu_\pi}(\boldsymbol{\eta}_{\boldsymbol{\theta}^*}, \boldsymbol{\eta}_{\bar{\boldsymbol{\theta}}_T})$ , whose proof is given in Appendix E.1.

**Lemma 5.1.** *For any  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$ , it holds that  $\mathcal{L}(\boldsymbol{\theta}) \leq 2K^{-1/2}(1-\gamma)^{-2}\lambda_{\min}^{-1/2} \|\bar{\mathbf{A}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$ .*

## 5.2 Exponential Stability Analysis

First, we introduce some notations. Letting  $\mathbf{e}_t := \mathbf{A}_t \boldsymbol{\theta}^* - \mathbf{b}_t$ , we denote by  $C_A$  (resp.  $C_e$ ) the almost sure upper bound for  $\max\{\|\mathbf{A}_t\|, \|\mathbf{A}_t - \bar{\mathbf{A}}\|\}$  (resp.  $\|\mathbf{e}_t\|$ ), and  $\boldsymbol{\Sigma}_e := \mathbb{E}[\mathbf{e}_t \mathbf{e}_t^\top]$  the covariance matrix of  $\mathbf{e}_t$ . The following lemma provides useful upper bounds, whose proof is given in Appendix E.2.

**Lemma 5.2.** *For any  $K \geq (1 - \gamma)^{-1}$ , it holds that*

$$C_A \leq 4, \quad C_e \leq 4(\|\boldsymbol{\theta}^*\| + \sqrt{K}(1 - \gamma)), \quad \text{tr}(\boldsymbol{\Sigma}_e) \leq 18(\|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi}^2 + K(1 - \gamma)^2).$$

Let  $\boldsymbol{\Gamma}_t^{(\alpha)} := \prod_{i=1}^t (\mathbf{I} - \alpha \mathbf{A}_i)$  for any  $\alpha > 0$  and  $t \in \mathbb{N}$ . The exponential stability of **Linear-CTD** is summarized in the following lemma, whose proof can be found in Appendix E.3.

**Lemma 5.3.** *For any  $p \geq 2$ , let  $a = (1 - \sqrt{\gamma})\lambda_{\min}/2$  and  $\alpha_{p,\infty} = (1 - \sqrt{\gamma})/(38p)$  ( $\alpha_{p,\infty} p \leq 1/2$ ). Then for any  $\alpha \in (0, \alpha_{p,\infty})$ ,  $\mathbf{u} \in \mathbb{R}^{dK}$  and  $t \in \mathbb{N}$ , it holds that  $\mathbb{E}^{1/p}[\|\boldsymbol{\Gamma}_t^{(\alpha)} \mathbf{u}\|^p] \leq (1 - \alpha a)^t \|\mathbf{u}\|$ .*

The following theorem states the  $L^2$  convergence of  $\|\bar{\mathbf{A}}(\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*)\|$  based on exponential stability arguments. For the general  $L^p$  convergence, please refer to Samsonov et al. [54, Theorem 2].

**Theorem 5.1.** [54, Theorem 1] *For any  $\alpha \in (0, \alpha_{2,\infty})$ , it holds that*

$$\mathbb{E}^{1/2}[\|\bar{\mathbf{A}}(\bar{\boldsymbol{\theta}}_T - \boldsymbol{\theta}^*)\|^2] \lesssim \frac{\sqrt{\text{tr}(\boldsymbol{\Sigma}_e)}}{\sqrt{T}} \left(1 + \frac{C_A \sqrt{\alpha}}{\sqrt{a}}\right) + \frac{\sqrt{\text{tr}(\boldsymbol{\Sigma}_e)}}{T \sqrt{\alpha a}} + \frac{(1 - \alpha a)^{T/2}}{T} \left(\frac{1}{\alpha} + \frac{C_A}{\sqrt{\alpha a}}\right) \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|.$$

Combining these lemmas with Theorem 5.1, we can immediately obtain Theorem 4.1.

**Remark 5** (Convergence of SSGD with PMF representation). *In Appendix E.5, we give counterparts of these lemmas and Theorem 4.1 for vanilla SSGD with PMF representation. The results imply that the step size in the algorithm should scale with  $(1 - \sqrt{\gamma})/K^2$  and the sample complexity grows with  $K$ . In Appendix F.2, we verify through numerical experiments that as  $K$  increases, to ensure convergence, the step size of the vanilla algorithm indeed needs to decay at a rate of  $K^{-2}$ . In contrast, the step size of our **Linear-CTD** does not need to be adjusted when  $K$  increases. Moreover, we find that **Linear-CTD** empirically consistently outperforms the vanilla algorithm under different  $K$ .*

## 5.3 Key Steps in the Proofs

Here we highlight some key steps in proving the above theoretical results.

**Bounding the Spectral Norm of Expectation of Kronecker Products.** In proving that the  $\mathcal{L}(\boldsymbol{\theta})$  can be upper-bounded by  $\|\bar{\mathbf{A}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$  (Lemma 5.1), as well as in verifying the exponential stability condition (Lemma 5.3), one of the most critical steps is to show

$$\left\| \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes \left( \boldsymbol{\Sigma}_{\phi}^{-\frac{1}{2}} \phi(s) \phi(s')^{\top} \boldsymbol{\Sigma}_{\phi}^{-\frac{1}{2}} \right) \right] \right\| \leq \sqrt{\gamma}, \quad \forall r \in [0, 1]. \quad (17)$$

By Lemma G.3, we have  $\|\mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1}\| \leq \sqrt{\gamma}$  for any  $r \in [0, 1]$ . In addition, one can check that  $\|\mathbb{E}_{s,s'}[\boldsymbol{\Sigma}_{\phi}^{-1/2} \phi(s) \phi(s')^{\top} \boldsymbol{\Sigma}_{\phi}^{-1/2}]\| \leq 1$ . One may speculate that the property  $\|\mathbf{B}_1 \otimes \mathbf{B}_2\| = \|\mathbf{B}_1\| \|\mathbf{B}_2\|$  (Lemma A.3) is enough to get the desired conclusion. However, the two matrices in the Kronecker product are not independent, preventing us from using this simple property to derive the conclusion. On the other hand, since we only have the upper bound  $\mathbb{E}_{s,s'}[\|\boldsymbol{\Sigma}_{\phi}^{-1/2} \phi(s) \phi(s')^{\top} \boldsymbol{\Sigma}_{\phi}^{-1/2}\|] \leq d$ , simply moving the expectation in Eqn. (17) outside the norm will lead to a loose  $d\sqrt{\gamma}$  bound. To resolve this problem, we leverage the fact that the second matrix is rank-1 and prove the following result. The proof can be found in the derivation following Eqn. (29).

**Lemma 5.4.** *For any random matrix  $\mathbf{Y}$  and random vectors  $\mathbf{x}, \mathbf{z}$ , suppose  $\|\mathbf{Y}\| \leq C_Y$  almost surely,  $\mathbb{E}[\mathbf{x}\mathbf{x}^{\top}] \leq C_x \mathbf{I}_{d_1}$  and  $\mathbb{E}[\mathbf{z}\mathbf{z}^{\top}] \leq C_z \mathbf{I}_{d_2}$  for some constants  $C_Y, C_x, C_z > 0$ . Then it holds that*

$$\|\mathbb{E}[\mathbf{Y} \otimes (\mathbf{x}\mathbf{z}^{\top})]\| \leq C_Y \sqrt{C_x C_z}.$$

**Remark 6** (Matrix Representation of Categorical Projected Bellman operator). *The matrix  $\mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1}$  also appears in Rowland et al. [51, Proposition B.2] as the matrix representation of the categorical projected Bellman operator  $\Pi_K \mathcal{T}^{\pi}$  of a specific one-state MDP. As a result,  $\|\mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1}\| \leq \sqrt{\gamma}$  because  $\Pi_K \mathcal{T}^{\pi}$  is a  $\sqrt{\gamma}$ -contraction in  $(\mathcal{P}, \ell_2)$ . Our Lemma G.3 provides a new analysis by directly analyzing the matrix.*

**Bounding the Norm of  $\mathbf{b}_t$ .** In proving Lemma 5.2, the most involved step is to upper-bound  $\|\mathbf{b}_t\|$  (Eqn. (16)). To this end, we need to upper-bound the following term:

$$\frac{1}{K+1} \left\| \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right) \right\|, \quad \forall r \in [0, 1]. \quad (18)$$

Term (18) is also related to the categorical projected Bellman operator  $\Pi_K \mathcal{T}^\pi$ . Specifically, let  $\nu = \frac{1}{K+1} \sum_{k=0}^K \delta_{x_k}$  be the discrete uniform distribution. One can show that Term (18) equals

$$\left\| C \left( p_{(b_{r,\gamma})_\#(\nu)} - p_\nu \right) \right\| = \iota_K^{-1/2} \ell_2(\Pi_K(b_{r,\gamma})_\#(\nu), \nu) \leq \iota_K^{-1/2} \ell_2((b_{r,\gamma})_\#(\nu), \nu) \leq 3\sqrt{K}(1-\gamma),$$

where we used the fact that  $\Pi_K$  is non-expansive and an upper bound for  $\ell_2((b_{r,\gamma})_\#(\nu), \nu)$  when  $K \geq (1-\gamma)^{-1}$  (Lemma G.4). The full proof can be found in the derivation following Eqn. (30).

## 6 Conclusions

In this paper, we have bridged a critical theoretical gap in distributional reinforcement learning by establishing the non-asymptotic sample complexity of distributional TD learning with linear function approximation. Specifically, we have proposed **Linear-CTD**, which is derived by solving the linear-categorical projected Bellman equation. By carefully analyzing the Bellman equation and using the exponential stability arguments, we have shown tight sample complexity bounds for the proposed algorithm. Our finite-sample rates match the state-of-the-art sample complexity bounds for conventional TD learning. These theoretical findings demonstrate that learning the full return distribution under linear function approximation can be statistically as easy as conventional TD learning for value function estimation. Our numerical experiments have provided empirical validation of our theoretical results. Finally, we have noted that it would be possible to improve the convergence rates by applying variance-reduction techniques or using polynomial-decaying step sizes, which we leave for future work.

## References

- [1] N. Bäuerle and J. Ott. Markov decision processes with average-value-at-risk criteria. *Mathematical Methods of Operations Research*, 74:361–379, 2011.
- [2] 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.
- [3] M. G. Bellemare, N. Le Roux, P. S. Castro, and S. Moitra. Distributional reinforcement learning with linear function approximation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2203–2211. PMLR, 2019.



- [4] M. G. Bellemare, W. Dabney, and M. Rowland. *Distributional Reinforcement Learning*. MIT Press, 2023. <http://www.distributional-rl.org>.
- [5] D. P. Bertsekas and J. N. Tsitsiklis. Neuro-dynamic programming: an overview. In *Proceedings of 1995 34th IEEE conference on decision and control*, volume 1, pages 560–564. IEEE, 1995.
- [6] J. Bhandari, D. Russo, and R. Singal. A finite time analysis of temporal difference learning with linear function approximation. In *Conference on learning theory*, pages 1691–1692. PMLR, 2018.
- [7] M. Böck and C. Heitzinger. Speedy categorical distributional reinforcement learning and complexity analysis. *SIAM Journal on Mathematics of Data Science*, 4(2):675–693, 2022. doi: 10.1137/20M1364436. URL <https://doi.org/10.1137/20M1364436>.
- [8] V. I. Bogachev. *Measure theory*, volume 1. Springer, 2007.
- [9] K. Chen. *Matrix preconditioning techniques and applications*. Cambridge University Press, 2005.
- [10] Z. Chen, S. T. Maguluri, S. Shakkottai, and K. Shanmugam. A lyapunov theory for finite-sample guarantees of markovian stochastic approximation. *Operations Research*, 72(4):1352–1367, 2024.
- [11] Y. Chow and M. Ghavamzadeh. Algorithms for cvar optimization in mdps. *Advances in neural information processing systems*, 27, 2014.
- [12] 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.
- [13] W. Dabney, M. Rowland, M. Bellemare, and R. Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [14] G. Dalal, B. Szörényi, G. Thoppe, and S. Mannor. Finite sample analyses for td (0) with function approximation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [15] Y. Duan and M. J. Wainwright. A finite-sample analysis of multi-step temporal difference estimates. In *Learning for Dynamics and Control Conference*, pages 612–624. PMLR, 2023.

- [16] A. Durmus, E. Moulines, A. Naumov, and S. Samsonov. Finite-time high-probability bounds for polyak–ruppert averaged iterates of linear stochastic approximation. *Mathematics of Operations Research*, 2024.
- [17] 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.
- [18] C. D. Godsil. Inverses of trees. *Combinatorica*, 5:33–39, 1985.
- [19] R. A. Horn and C. R. Johnson. *Topics in matrix analysis*. Cambridge university press, 1994.
- [20] D. Huo, Y. Chen, and Q. Xie. Bias and extrapolation in markovian linear stochastic approximation with constant stepsizes. In *Abstract Proceedings of the 2023 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, pages 81–82, 2023.
- [21] T. Kastner, M. A. Erdogdu, and A.-m. Farahmand. Distributional model equivalence for risk-sensitive reinforcement learning. *Advances in Neural Information Processing Systems*, 36: 56531–56552, 2023.
- [22] T. L. Lai. Stochastic approximation. *The Annals of Statistics*, 31(2):391–406, 2003.
- [23] C. Lakshminarayanan and C. Szepesvari. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *International conference on artificial intelligence and statistics*, pages 1347–1355. PMLR, 2018.
- [24] P. Lancaster and H. K. Farahat. Norms on direct sums and tensor products. *mathematics of computation*, 26(118):401–414, 1972.
- [25] G. Li, C. Cai, Y. Chen, Y. Wei, and Y. Chi. Is q-learning minimax optimal? a tight sample complexity analysis. *Operations Research*, 72(1):222–236, 2024.
- [26] G. Li, W. Wu, Y. Chi, C. Ma, A. Rinaldo, and Y. Wei. High-probability sample complexities for policy evaluation with linear function approximation. *IEEE Transactions on Information Theory*, 2024.

- [27] T. Li, G. Lan, and A. Pananjady. Accelerated and instance-optimal policy evaluation with linear function approximation. *SIAM Journal on Mathematics of Data Science*, 5(1):174–200, 2023.
- [28] X. Li, J. Liang, and Z. Zhang. Online statistical inference for nonlinear stochastic approximation with markovian data. *arXiv preprint arXiv:2302.07690*, 2023.
- [29] X.-L. Li. Preconditioned stochastic gradient descent. *IEEE transactions on neural networks and learning systems*, 29(5):1454–1466, 2017.
- [30] S. H. Lim and I. Malik. Distributional reinforcement learning for risk-sensitive policies. *Advances in Neural Information Processing Systems*, 35:30977–30989, 2022.
- [31] Y. Luo, G. Liu, H. Duan, O. Schulte, and P. Poupart. Distributional reinforcement learning with monotonic splines. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=C8Ltz08PtBp>.
- [32] C. Lyle, M. G. Bellemare, and P. S. Castro. A comparative analysis of expected and distributional reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4504–4511, 2019.
- [33] A. Marthe, A. Garivier, and C. Vernade. Beyond average return in markov decision processes. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=mgNu8nDFwa>.
- [34] S. P. Meyn and R. L. Tweedie. *Markov chains and stochastic stability*. Springer Science & Business Media, 2012.
- [35] M. Moghimi and H. Ku. Beyond cvar: Leveraging static spectral risk measures for enhanced decision-making in distributional reinforcement learning. *arXiv preprint arXiv:2501.02087*, 2025.
- [36] 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.
- [37] W. Mou, C. J. Li, M. J. Wainwright, P. L. Bartlett, and M. I. Jordan. On linear stochastic approximation: Fine-grained polyak-ruppert and non-asymptotic concentration. In *Conference on Learning Theory*, pages 2947–2997. PMLR, 2020.

- [38] W. Mou, A. Pananjady, M. Wainwright, and P. Bartlett. Optimal and instance-dependent guarantees for markovian linear stochastic approximation. In *Conference on Learning Theory*, pages 2060–2061. PMLR, 2022.
- [39] T. Nguyen-Tang, S. Gupta, and S. Venkatesh. Distributional reinforcement learning via moment matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9144–9152, 2021.
- [40] E. Noorani, C. N. Mavridis, and J. S. Baras. Exponential td learning: A risk-sensitive actor-critic reinforcement learning algorithm. In *2023 American Control Conference (ACC)*, pages 4104–4109. IEEE, 2023.
- [41] G. Patil, L. Prashanth, D. Nagaraj, and D. Precup. Finite time analysis of temporal difference learning with linear function approximation: Tail averaging and regularisation. In *International Conference on Artificial Intelligence and Statistics*, pages 5438–5448. PMLR, 2023.
- [42] D. Paulin. Concentration inequalities for markov chains by marton couplings and spectral methods. *Electronic journal of probability*, 20:79, 2015.
- [43] Y. Peng, L. Zhang, and Z. Zhang. Statistical efficiency of distributional temporal difference learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=eWUM5hRYgH>.
- [44] B. Á. Pires, M. Rowland, D. Borsa, Z. D. Guo, K. Khetarpal, A. Barreto, D. Abel, R. Munos, and W. Dabney. Optimizing return distributions with distributional dynamic programming. *arXiv preprint arXiv:2501.13028*, 2025.
- [45] B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- [46] C. Qu, S. Mannor, and H. Xu. Nonlinear distributional gradient temporal-difference learning. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5251–5260. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/qu19b.html>.
- [47] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.

- [48] 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.
- [49] M. Rowland, Y. Tang, C. Lyle, R. Munos, M. G. Bellemare, and W. Dabney. The statistical benefits of quantile temporal-difference learning for value estimation. In *International Conference on Machine Learning*, pages 29210–29231. PMLR, 2023.
- [50] 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. *Journal of Machine Learning Research*, 25:1–47, 2024.
- [51] M. Rowland, L. K. Wenliang, R. Munos, C. Lyle, Y. Tang, and W. Dabney. Near-minimax-optimal distributional reinforcement learning with a generative model. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=JXKbf1d4ib>.
- [52] D. Ruppert. Efficient estimations from a slowly convergent robbins-monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- [53] S. Samsonov, E. Moulines, Q.-M. Shao, Z.-S. Zhang, and A. Naumov. Gaussian approximation and multiplier bootstrap for polyak-ruppert averaged linear stochastic approximation with applications to TD learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=S0Ci1AsJL5>.
- [54] S. Samsonov, D. Tiapkin, A. Naumov, and E. Moulines. Improved high-probability bounds for the temporal difference learning algorithm via exponential stability. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 4511–4547. PMLR, 2024.
- [55] D. Serre. *Matrices: Theory and Applications*. Graduate texts in mathematics. Springer, 2002. ISBN 9780387954608. URL <https://books.google.com.hk/books?id=RDnUIFYgrUC>.
- [56] R. Srikant and L. Ying. Finite-time error bounds for linear stochastic approximation and td learning. In *Conference on Learning Theory*, pages 2803–2830. PMLR, 2019.
- [57] K. Sun, Y. Zhao, W. Liu, B. Jiang, and L. Kong. Distributional reinforcement learning with regularized wasserstein loss. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=CiEynTpF28>.

- [58] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3: 9–44, 1988.
- [59] Y. Tang, R. Munos, M. Rowland, B. Avila Pires, W. Dabney, and M. Bellemare. The nature of temporal difference errors in multi-step distributional reinforcement learning. *Advances in Neural Information Processing Systems*, 35:30265–30276, 2022.
- [60] Y. Tang, M. Rowland, R. Munos, B. Á. Pires, and W. Dabney. Off-policy distributional q ( $\lambda$ ): Distributional rl without importance sampling. *arXiv preprint arXiv:2402.05766*, 2024.
- [61] J. Tsitsiklis and B. Van Roy. Analysis of temporal-difference learning with function approximation. *Advances in neural information processing systems*, 9, 1996.
- [62] 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.
- [63] K. Wang, K. Zhou, R. Wu, N. Kallus, and W. Sun. The benefits of being distributional: Small-loss bounds for reinforcement learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 2275–2312. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/06fc38f5c21ae66ef955e28b7a78ece5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/06fc38f5c21ae66ef955e28b7a78ece5-Paper-Conference.pdf).
- [64] K. Wang, O. Oertell, A. Agarwal, N. Kallus, and W. Sun. More benefits of being distributional: Second-order bounds for reinforcement learning. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=kZBCFQe1Ej>.
- [65] L. K. Wenliang, G. Deletang, M. Aitchison, M. Hutter, A. Ruoss, A. Gretton, and M. Rowland. Distributional bellman operators over mean embeddings. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=j2pLfsBm4J>.
- [66] H. Wiltzer, J. Farebrother, A. Gretton, and M. Rowland. Foundations of multivariate distributional reinforcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=aq3I5B6GLG>.

- [67] H. Wiltzer, J. Farebrother, A. Gretton, Y. Tang, A. Barreto, W. Dabney, M. G. Bellemare, and M. Rowland. A distributional analogue to the successor representation. In *International Conference on Machine Learning (ICML)*, 2024.
- [68] R. Wu, M. Uehara, and W. Sun. Distributional offline policy evaluation with predictive error guarantees. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 37685–37712. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/wu23s.html>.
- [69] W. Wu, G. Li, Y. Wei, and A. Rinaldo. Statistical Inference for Temporal Difference Learning with Linear Function Approximation. *arXiv preprint arXiv:2410.16106*, 2024.
- [70] D. Yang, L. Zhao, Z. Lin, T. Qin, J. Bian, and T.-Y. Liu. Fully parameterized quantile function for distributional reinforcement learning. *Advances in neural information processing systems*, 32, 2019.
- [71] Y. Yue, Z. Wang, and M. Zhou. Implicit distributional reinforcement learning. *Advances in Neural Information Processing Systems*, 33:7135–7147, 2020.
- [72] H. Zhang and F. Ding. On the kronecker products and their applications. *Journal of Applied Mathematics*, 2013(1):296185, 2013.
- [73] L. Zhang, Y. Peng, J. Liang, W. Yang, and Z. Zhang. Estimation and Inference in Distributional Reinforcement Learning. *The Annals of Statistics*, 2025.
- [74] F. Zhou, Z. Zhu, Q. Kuang, and L. Zhang. Non-decreasing quantile function network with efficient exploration for distributional reinforcement learning. In Z.-H. Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 3455–3461. International Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/476. URL <https://doi.org/10.24963/ijcai.2021/476>. Main Track.

## A Kronecker Product

In this section, we will introduce some properties of Kronecker product used in our paper. See Zhang and Ding [72] for a detailed treatment of Kronecker product.

For any matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{p \times q}$ , the Kronecker product  $\mathbf{A} \otimes \mathbf{B}$  is a matrix in  $\mathbb{R}^{mp \times nq}$ , defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}.$$

**Lemma A.1.** *The Kronecker product is bilinear and associative. Furthermore, for any matrices  $\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3, \mathbf{B}_4$  such that  $\mathbf{B}_1\mathbf{B}_3, \mathbf{B}_2\mathbf{B}_4$  can be defined, it holds that  $(\mathbf{B}_1 \otimes \mathbf{B}_2)(\mathbf{B}_3 \otimes \mathbf{B}_4) = (\mathbf{B}_1\mathbf{B}_3) \otimes (\mathbf{B}_2\mathbf{B}_4)$  (mixed-product property).*

*Proof.* See Basic properties and Theorem 3 in Zhang and Ding [72].  $\square$

**Lemma A.2.** *For any matrices  $\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3$  such that  $\mathbf{B}_1\mathbf{B}_2\mathbf{B}_3$  can be defined, it holds that  $\text{vec}(\mathbf{B}_1\mathbf{B}_2\mathbf{B}_3) = (\mathbf{B}_3^\top \otimes \mathbf{B}_1) \text{vec}(\mathbf{B}_2)$ .*

*Proof.* See Lemma 4.3.1 in Horn and Johnson [19].  $\square$

**Lemma A.3.** *For any matrices  $\mathbf{B}_1$  and  $\mathbf{B}_2$ , it holds that  $\|\mathbf{B}_1 \otimes \mathbf{B}_2\| = \|\mathbf{B}_1\| \|\mathbf{B}_2\|$ ,  $(\mathbf{B}_1 \otimes \mathbf{B}_2)^\top = \mathbf{B}_1^\top \otimes \mathbf{B}_2^\top$ . Furthermore, if  $\mathbf{B}_1$  and  $\mathbf{B}_2$  are invertible/orthogonal/diagonal/symmetric/normal,  $\mathbf{B}_1 \otimes \mathbf{B}_2$  is also invertible/orthogonal/diagonal/symmetric/normal and  $(\mathbf{B}_1 \otimes \mathbf{B}_2)^{-1} = \mathbf{B}_1^{-1} \otimes \mathbf{B}_2^{-1}$ .*

*Proof.* See Basic properties, Theorem 5 and Theorem 7 in Zhang and Ding [72].  $\square$

**Lemma A.4.** *For any  $K, d \in \mathbb{N}$  and PSD matrices  $\mathbf{B}_1, \mathbf{B}_3 \in \mathbb{R}^{K \times K}, \mathbf{B}_2, \mathbf{B}_4 \in \mathbb{R}^{d \times d}$  with  $\mathbf{B}_1 \preceq \mathbf{B}_3$  and  $\mathbf{B}_2 \preceq \mathbf{B}_4$ , it holds that  $\mathbf{B}_1 \otimes \mathbf{B}_2, \mathbf{B}_3 \otimes \mathbf{B}_4$  are also PSD matrices, furthermore,  $\mathbf{B}_1 \otimes \mathbf{B}_2 \preceq \mathbf{B}_3 \otimes \mathbf{B}_4$ .*

*Proof.* Consider the spectral decomposition  $\mathbf{B}_i = \mathbf{Q}_i \mathbf{D}_i \mathbf{Q}_i^\top$ , for any  $i \in [4]$ , by Lemma A.1 and Lemma A.3, we have

$$(\mathbf{B}_1 \otimes \mathbf{B}_2) = (\mathbf{Q}_1 \otimes \mathbf{Q}_2) (\mathbf{D}_1 \otimes \mathbf{D}_2) (\mathbf{Q}_1 \otimes \mathbf{Q}_2)^\top$$

and

$$(\mathbf{B}_3 \otimes \mathbf{B}_4) = (\mathbf{Q}_3 \otimes \mathbf{Q}_4) (\mathbf{D}_3 \otimes \mathbf{D}_4) (\mathbf{Q}_3 \otimes \mathbf{Q}_4)^\top$$



are also spectral decomposition of  $(\mathbf{B}_1 \otimes \mathbf{B}_2)$  and  $(\mathbf{B}_3 \otimes \mathbf{B}_4)$  respectively. It is easy to see that they are PSD. Furthermore,

$$\begin{aligned} (\mathbf{B}_3 \otimes \mathbf{B}_4) - (\mathbf{B}_1 \otimes \mathbf{B}_2) &= [(\mathbf{B}_3 \otimes \mathbf{B}_4) - (\mathbf{B}_3 \otimes \mathbf{B}_2)] + [(\mathbf{B}_3 \otimes \mathbf{B}_2) - (\mathbf{B}_1 \otimes \mathbf{B}_2)] \\ &= [\mathbf{B}_3 \otimes (\mathbf{B}_4 - \mathbf{B}_2)] + [(\mathbf{B}_3 - \mathbf{B}_1) \otimes \mathbf{B}_2] \\ &\geq \mathbf{0}. \end{aligned}$$

□

**Lemma A.5.** *For any  $K, d, d_1, d_2 \in \mathbb{N}$ , vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$  and matrices  $\mathbf{B}_1 \in \mathbb{R}^{K \times d_1}$ ,  $\mathbf{B}_2 \in \mathbb{R}^{d_2 \times K}$ ,  $\mathbf{B}_3 \in \mathbb{R}^{K \times K}$ , it holds that*

$$\begin{aligned} (\mathbf{I}_K \otimes \mathbf{u}) \mathbf{B}_1 &= \mathbf{B}_1 \otimes \mathbf{u}, \\ \mathbf{B}_2 (\mathbf{I}_K \otimes \mathbf{v})^\top &= \mathbf{B}_2 \otimes \mathbf{v}^\top, \\ (\mathbf{I}_K \otimes \mathbf{u}) \mathbf{B}_3 (\mathbf{I}_K \otimes \mathbf{v})^\top &= \mathbf{B}_3 \otimes (\mathbf{u} \mathbf{v}^\top). \end{aligned}$$

Furthermore, for any matrix  $\mathbf{B}_4 \in \mathbb{R}^{d_1 \times d_2}$ , we have

$$(\mathbf{B}_1 \otimes \mathbf{u}) \mathbf{B}_4 = (\mathbf{B}_1 \mathbf{B}_4) \otimes \mathbf{u}.$$

*Proof.* Let  $\mathbf{u} = (u_i)_{i=1}^d$   $\mathbf{B}_1 = (b_{ij})_{i,j=1}^K$ , then

$$\begin{aligned}
(\mathbf{I}_K \otimes \mathbf{u}) \mathbf{B}_1 &= \begin{bmatrix} \mathbf{u} & \mathbf{0}_d & \cdots & \mathbf{0}_d & \mathbf{0}_d \\ \mathbf{0}_d & \mathbf{u} & \cdots & \mathbf{0}_d & \mathbf{0}_d \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0}_d & \mathbf{0}_d & \cdots & \mathbf{u} & \mathbf{0}_d \\ \mathbf{0}_d & \mathbf{0}_d & \cdots & \mathbf{0}_d & \mathbf{u} \end{bmatrix} \begin{bmatrix} b_{11} & \cdots & b_{1K} \\ \vdots & \ddots & \vdots \\ b_{K1} & \cdots & b_{KK} \end{bmatrix} \\
&= \begin{bmatrix} b_{11}u_1 & \cdots & b_{1K}u_1 \\ \vdots & \ddots & \vdots \\ b_{11}u_d & \cdots & b_{1K}u_d \\ \vdots & \ddots & \vdots \\ b_{K1}u_1 & \cdots & b_{KK}u_1 \\ \vdots & \ddots & \vdots \\ b_{K1}u_d & \cdots & b_{KK}u_d \end{bmatrix} \\
&= \begin{bmatrix} b_{11}\mathbf{u} & \cdots & b_{1K}\mathbf{u} \\ \vdots & \ddots & \vdots \\ b_{K1}\mathbf{u} & \cdots & b_{KK}\mathbf{u} \end{bmatrix} \\
&= \mathbf{B}_1 \otimes \mathbf{u}.
\end{aligned}$$

Hence

$$\mathbf{B}_2 (\mathbf{I}_K \otimes \mathbf{v})^\top = [(\mathbf{I}_K \otimes \mathbf{v}) \otimes \mathbf{B}_2^\top]^\top = [\mathbf{B}_2^\top \otimes \mathbf{v}]^\top = \mathbf{B}_2 \otimes \mathbf{v}^\top.$$

And in the same way,

$$(\mathbf{I}_K \otimes \mathbf{u}) \mathbf{B}_3 (\mathbf{I}_K \otimes \mathbf{v})^\top = (\mathbf{B}_3 \otimes \mathbf{u}) \otimes \mathbf{v}^\top = \mathbf{B}_3 \otimes (\mathbf{u} \otimes \mathbf{v}^\top) = \mathbf{B}_3 \otimes (\mathbf{u} \mathbf{v}^\top).$$

Furthermore,

$$(\mathbf{B}_1 \otimes \mathbf{u}) \mathbf{B}_4 = [(\mathbf{I}_K \otimes \mathbf{u}) \mathbf{B}_1] \mathbf{B}_4 = (\mathbf{I}_K \otimes \mathbf{u}) (\mathbf{B}_1 \mathbf{B}_4) = (\mathbf{B}_1 \mathbf{B}_4) \otimes \mathbf{u}.$$

□

## B Related Work

**Distributional Reinforcement Learning.** Distributional TD learning was first proposed in Bellemare et al. [2]. Following the distributional perspective in Bellemare et al. [2], Qu et al. [46] proposed a distributional version of the gradient TD learning algorithm, Tang et al. [59] proposed a distributional version of multi-step TD learning, Tang et al. [60] proposed a distributional version of off-policy  $Q(\lambda)$  and  $TD(\lambda)$  algorithms, and Wu et al. [68] proposed a distributional version of fitted  $Q$  evaluation to solve the distributional offline policy evaluation problem. Wiltzer et al. [67] proposed an approach for evaluating the return distributions for all policies simultaneously when the reward is deterministic or in the finite-horizon setting. Wiltzer et al. [66] studied distributional policy evaluation in the multivariate reward setting and proposed corresponding TD learning algorithms. Beyond the tabular setting, Bellemare et al. [3], Lyle et al. [32], Bellemare et al. [4] proposed various distributional TD learning algorithms with linear function approximation under different parametrizations.

A series of recent studies have focused on the theoretical properties of distributional TD learning. Rowland et al. [48], Böck and Heitzinger [7], Zhang et al. [73], Rowland et al. [50, 51], Peng et al. [43] analyzed the asymptotic and non-asymptotic convergence of distributional TD learning (or its model-based variants) in the tabular setting. Among these works, Rowland et al. [51], Peng et al. [43] established that in the tabular setting, learning the full return distribution is statistically as easy as learning its expectation in the model-based and model-free settings, respectively. And Bellemare et al. [3] provided an asymptotic convergence result for categorical TD learning with linear function approximation.

Beyond the problem of distributional policy evaluation, Rowland et al. [49], Wang et al. [63, 64] showed that theoretically the classic value-based reinforcement learning could benefit from distributional reinforcement learning. Bäuerle and Ott [1], Chow and Ghavamzadeh [11], Marthe et al. [33], Noorani et al. [40], Moghimi and Ku [35], Pires et al. [44] considered optimizing statistical functionals of the return, and proposed algorithms to solve this harder problem.

**Stochastic Approximation.** Our **Linear-CTD** falls into the category of LSA. The classic TD learning, as one of the most classic LSA problems, has been extensively studied [5, 61, 6, 14, 41, 15, 25, 26, 53, 69]. Among these works, Li et al. [26], Samsonov et al. [54] provided the tightest bounds for **Linear-TD** with constant step sizes, which is also considered in our paper. While Wu et al. [69] established the tightest bounds for **Linear-TD** with polynomial-decaying step sizes.

For general stochastic approximation problems, extensive works [23, 56, 37, 38, 20, 28, 16, 54, 10] have provided solid theoretical understandings.

## C Omitted Results and Proofs in Section 2

### C.1 Linear Projected Bellman Equation

It is worth noting that,  $\Pi_\phi^\pi: (\mathbb{R}^S, \|\cdot\|_{\mu_\pi}) \rightarrow (\mathcal{V}_\phi, \|\cdot\|_{\mu_\pi})$  is an orthogonal projection.

We aim to derive Eqn. (4). It is easy to check that, for any  $\mathbf{V} \in \mathbb{R}^S$ ,  $\Pi_\phi^\pi \mathbf{V}$  is uniquely give by  $\mathbf{V}_{\tilde{\psi}}$  where

$$\tilde{\psi} = \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} [\phi(s) V(s)].$$

Hence, by the definition of Bellman operator (Eqn. (1)),  $\psi^\star$  is the unique solution to the following system of linear equations for  $\psi \in \mathbb{R}^d$

$$\begin{aligned} \psi &= \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} [\phi(s) [T^\pi \mathbf{V}_\psi](s)] \\ &= \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} [\phi(s) (\mathbb{E}[r_0 | s_0 = s] + \gamma \mathbb{E}[\phi(s_1)^\top | s_0 = s] \psi)] \\ &= \Sigma_\phi^{-1} \mathbb{E}_{s, s'} [\phi(s) \phi(s')^\top] \psi + \Sigma_\phi^{-1} \mathbb{E}_{s, r} [\phi(s) r], \end{aligned}$$

or equivalently,

$$(\Sigma_\phi - \gamma \mathbb{E}_{s, s'} [\phi(s) \phi(s')^\top]) \psi = \mathbb{E}_{s, r} [\phi(s) r].$$

### C.2 Convergence Results for Linear TD Learning

It is worthy noting that, **Linear**-TD is equivalent to the stochastic semi-gradient descent (SSGD) update.

In **Linear**-TD, our goal is to find a good estimator  $\hat{\psi}$  such that  $\|\mathbf{V}_{\hat{\psi}} - \mathbf{V}_{\psi^\star}\|_{\mu_\pi} = \|\hat{\psi} - \psi^\star\|_{\Sigma_\phi} \leq \varepsilon$ . [54] considered the Polyak-Ruppert tail averaging  $\bar{\psi}_T := (T/2 + 1)^{-1} \sum_{t=T/2}^T \psi_t$ , and showed that in the generative model setting with constant step size  $\alpha \simeq (1 - \gamma)\lambda_{\min}$ ,

$$T = \tilde{O} \left( \frac{\|\psi^\star\|_{\Sigma_\phi}^2 + 1}{(1 - \gamma)^2 \lambda_{\min}} \left( \frac{1}{\varepsilon^2} + \frac{1}{\lambda_{\min}} \right) \right)$$

is sufficient to guarantee that  $\|\mathbf{V}_{\bar{\psi}_T} - \mathbf{V}_{\psi^\star}\|_{\mu_\pi} \leq \varepsilon$ . They also provided sample complexity bounds when taking the instance-independent (*i.e.*, not dependent on unknown quantity) step size, and in the Markovian setting.

### C.3 Categorical Parametrization is an Isometry

**Proposition C.1.** *The affine space  $(\mathcal{P}_K^{\text{sign}}, \ell_2)$  is isometric with  $(\mathbb{R}^K, \sqrt{\iota_K} \|\cdot\|_{\mathbf{C}^\top \mathbf{C}})$ , in the sense that, for any  $\nu_{\mathbf{p}_1}, \nu_{\mathbf{p}_2} \in \mathcal{P}_K^{\text{sign}}$ , it holds that  $\ell_2^2(\nu_{\mathbf{p}_1}, \nu_{\mathbf{p}_2}) = \iota_K \|\mathbf{p}_1 - \mathbf{p}_2\|_{\mathbf{C}^\top \mathbf{C}}^2$ , where  $\mathbf{C}$  is defined in Eqn. (11).*

*Proof.*

$$\begin{aligned} \ell_2^2(\nu_{\mathbf{p}_1}, \nu_{\mathbf{p}_2}) &= \int_0^{(1-\gamma)^{-1}} (F_{\nu_{\mathbf{p}_1}}(x) - F_{\nu_{\mathbf{p}_2}}(x))^2 dx \\ &= \iota_K \sum_{k=0}^{K-1} (F_{\nu_{\mathbf{p}_1}}(x_k) - F_{\nu_{\mathbf{p}_2}}(x_k))^2 \\ &= \iota_K \|\mathbf{C}(\mathbf{p}_1 - \mathbf{p}_2)\|^2 \\ &= \iota_K \|\mathbf{p}_1 - \mathbf{p}_2\|_{\mathbf{C}^\top \mathbf{C}}^2. \end{aligned}$$

□

### C.4 Categorical Projection Operator is Orthogonal Projection

**Proposition C.2.** [4, Lemma 9.17] *For any  $\nu \in \mathcal{P}^{\text{sign}}$  and  $\nu_{\mathbf{p}} \in \mathcal{P}_K^{\text{sign}}$ , it holds that*

$$\ell_2^2(\nu, \nu_{\mathbf{p}}) = \ell_2^2(\nu, \mathbf{\Pi}_K \nu) + \ell_2^2(\mathbf{\Pi}_K \nu, \nu_{\mathbf{p}}).$$

### C.5 Categorical Projected Bellman Operator

The following lemma characterizing  $\mathbf{\Pi}_K \mathcal{T}^\pi$  is useful for both practice and theoretical analysis.

**Proposition C.3.** *For any  $\boldsymbol{\eta} \in (\mathcal{P}^{\text{sign}})^{\mathcal{S}}$  and  $s \in \mathcal{S}$ , it holds that*

$$\begin{aligned} \mathbf{p}_{\mathcal{T}^\pi \boldsymbol{\eta}}(s) &= \mathbb{E} \left[ \mathbf{g}_K(r_0) + (\mathbf{G}(r_0) - \mathbf{1}_K^\top \otimes \mathbf{g}_K(r_0)) \mathbf{p}_{\boldsymbol{\eta}}(s_1) \middle| s_0 = s \right] \\ &= \mathbb{E} \left[ \tilde{\mathbf{G}}(r_0) \left( \mathbf{p}_{\boldsymbol{\eta}}(s_1) - \frac{1}{K+1} \mathbf{1}_K \right) \middle| s_0 = s \right] + \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ \mathbf{g}_j(r_0) \middle| s_0 = s \right]. \end{aligned}$$

And in the same way, for any  $r \in [0, 1]$  and  $s' \in \mathcal{S}$ , it holds that

$$\mathbf{p}_{(b_{r,\gamma})_{\#} \boldsymbol{\eta}(s')} = \tilde{\mathbf{G}}(r) \left( \mathbf{p}_{\boldsymbol{\eta}}(s') - \frac{1}{K+1} \mathbf{1}_K \right) + \frac{1}{K+1} \sum_{j=0}^K \mathbf{g}_j(r),$$

where  $\tilde{\mathbf{G}}$  and  $\mathbf{g}$  is defined in Theorem 3.1.

This proposition is a special case of Proposition D.3, whose proof can be found in Appendix D.4.

## D Omitted Results and Proofs in Section 3

### D.1 Linear-Categorical Parametrization is an Isometry

**Proposition D.1.** *The affine space  $(\mathcal{P}_{\phi,K}^{\text{sign}}, \ell_{2,\mu_\pi})$  is isometric with  $(\mathbb{R}^{dK}, \sqrt{\iota_K} \|\cdot\|_{\mathbf{I}_K \otimes \Sigma_\phi})$ , in the sense that, for any  $\boldsymbol{\eta}_{\theta_1}, \boldsymbol{\eta}_{\theta_2} \in \mathcal{P}_{\phi,K}^{\text{sign}}$ , it holds that  $\ell_{2,\mu_\pi}^2(\boldsymbol{\eta}_{\theta_1}, \boldsymbol{\eta}_{\theta_2}) = \iota_K \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_{\mathbf{I}_K \otimes \Sigma_\phi}^2$ .*

*Proof.* For any  $\boldsymbol{\eta}_\theta \in \mathcal{P}_{\phi,K}^{\text{sign}}$ , we denote  $\mathbf{F}_\theta(s) = (F_k(s; \boldsymbol{\theta}))_{k=0}^{K-1} \in \mathbb{R}^K$ , then it holds that

$$\begin{aligned} \ell_{2,\mu_\pi}^2(\boldsymbol{\eta}_{\theta_1}, \boldsymbol{\eta}_{\theta_2}) &= \iota_K \mathbb{E}_{s \sim \mu_\pi} [\|\mathbf{F}_{\theta_1}(s) - \mathbf{F}_{\theta_2}(s)\|^2] \\ &= \iota_K \text{tr} \left( \Sigma_\phi^{\frac{1}{2}} (\boldsymbol{\Theta}_1 - \boldsymbol{\Theta}_2) (\boldsymbol{\Theta}_1 - \boldsymbol{\Theta}_2)^\top \Sigma_\phi^{\frac{1}{2}} \right) \\ &= \iota_K \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_{\mathbf{I}_K \otimes \Sigma_\phi}^2. \end{aligned}$$

□

### D.2 Linear-Categorical Projection Operator

Proposition 3.1 is an immediate corollary of the following lemma. For any  $\nu \in \mathcal{P}_K^{\text{sign}}$ , we define  $\mathbf{F}_\nu = (F_k(\nu))_{k=0}^{K-1} = (\nu([0, x_k]))_{k=0}^{K-1} \in \mathbb{R}^K$ , and for any  $\boldsymbol{\eta} \in (\mathcal{P}^{\text{sign}})^\mathcal{S}$ , we define  $\mathbf{p}_\eta(s) = \mathbf{p}_{\Pi_K \eta(s)}$  and  $\mathbf{F}_\eta(s) = \mathbf{F}_{\Pi_K \eta(s)}$ .

**Lemma D.1.** *For any  $\boldsymbol{\eta} \in (\mathcal{P}^{\text{sign}})^\mathcal{S}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$  and  $s \in \mathcal{S}$ , it holds that*

$$\begin{aligned} \nabla_{\boldsymbol{\Theta}} \ell_2^2(\boldsymbol{\eta}_\theta(s), \boldsymbol{\eta}(s)) &= 2\iota_K \phi(s) (\mathbf{F}_\theta(s) - \mathbf{F}_\eta(s))^\top \\ &= 2\iota_K \phi(s) \left[ \phi(s)^\top \boldsymbol{\Theta} + \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_\eta(s) \right)^\top \mathbf{C}^\top \right]. \end{aligned} \tag{19}$$

Furthermore, it holds that

$$\begin{aligned} \nabla_{\boldsymbol{\Theta}} \ell_{2,\mu_\pi}^2(\boldsymbol{\eta}_\theta, \boldsymbol{\eta}) &= \mathbb{E}_{s \sim \mu_\pi} [\nabla_{\boldsymbol{\Theta}} \ell_2^2(\boldsymbol{\eta}_\theta(s), \boldsymbol{\eta}(s))] \\ &= 2\iota_K \left[ \Sigma_\phi \boldsymbol{\Theta} + \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_\eta(s) \right)^\top \right] \mathbf{C}^\top \right]. \end{aligned}$$

*Proof.* According to Proposition C.2, one has

$$\ell_2^2(\boldsymbol{\eta}_\theta(s), \boldsymbol{\eta}(s)) = \ell_2^2(\boldsymbol{\eta}_\theta(s), \boldsymbol{\Pi}_K \boldsymbol{\eta}(s)) + \ell_2^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}(s), \boldsymbol{\eta}(s)).$$

Hence,

$$\begin{aligned}
\nabla_{\boldsymbol{\theta}} \ell_2^2(\eta_{\boldsymbol{\theta}}(s), \eta(s)) &= \nabla_{\boldsymbol{\theta}} \ell_2^2(\eta_{\boldsymbol{\theta}}(s), \boldsymbol{\Pi}_K \eta(s)) \\
&= \iota_K \nabla_{\boldsymbol{\theta}} \|\mathbf{F}_{\boldsymbol{\theta}}(s) - \mathbf{F}_{\boldsymbol{\eta}}(s)\|^2 \\
&= 2\iota_K (\mathbf{I}_K \otimes \phi(s)) (\mathbf{F}_{\boldsymbol{\theta}}(s) - \mathbf{F}_{\boldsymbol{\eta}}(s)) \\
&= 2\iota_K (\mathbf{I}_K \otimes \phi(s)) \left( (\mathbf{I}_K \otimes \phi(s))^\top \boldsymbol{\theta} + \mathbf{C} \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\boldsymbol{\eta}}(s) \right) \right) \\
&= 2\iota_K \left\{ [\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)] \boldsymbol{\theta} + \left[ \left( \mathbf{C} \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\boldsymbol{\eta}}(s) \right) \right) \otimes \phi(s) \right] \right\}.
\end{aligned}$$

We also have the following matrix representation:

$$\begin{aligned}
\nabla_{\boldsymbol{\Theta}} \ell_2^2(\eta_{\boldsymbol{\theta}}(s), \eta(s)) &= 2\iota_K \phi(s) (\mathbf{F}_{\boldsymbol{\theta}}(s) - \mathbf{F}_{\boldsymbol{\eta}}(s))^\top \\
&= 2\iota_K \phi(s) \left[ \phi(s)^\top \boldsymbol{\Theta} + \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\boldsymbol{\eta}}(s) \right)^\top \mathbf{C}^\top \right].
\end{aligned}$$

□

**Proposition D.2.** For any  $\boldsymbol{\eta} \in (\mathcal{P}^{\text{sign}})^{\mathcal{S}}$  and  $\boldsymbol{\eta}_{\boldsymbol{\theta}} \in \mathcal{P}_{\phi, K}^{\text{sign}}$ , it holds that

$$\ell_{2, \mu_\pi}^2(\boldsymbol{\eta}, \boldsymbol{\eta}_{\boldsymbol{\theta}}) = \ell_{2, \mu_\pi}^2(\boldsymbol{\eta}, \boldsymbol{\Pi}_K \boldsymbol{\eta}) + \ell_{2, \mu_\pi}^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}, \boldsymbol{\Pi}_{\phi, K}^\pi \boldsymbol{\eta}) + \ell_{2, \mu_\pi}^2(\boldsymbol{\Pi}_{\phi, K}^\pi \boldsymbol{\eta}, \boldsymbol{\eta}_{\boldsymbol{\theta}}).$$

The proof is straightforward and almost the same as that of Proposition C.2 if we utilize the affine structure.

### D.3 Linear-Categorical Projected Bellman Equation

To derive the result, the following proposition characterizing  $\boldsymbol{\Pi}_K \mathcal{T}^\pi \boldsymbol{\eta}_{\boldsymbol{\theta}}$  is useful, whose proof can be found in Appendix D.4.

**Proposition D.3.** For any  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$  and  $s \in \mathcal{S}$ , we abbreviate  $\mathbf{p}_{\mathcal{T}^\pi \boldsymbol{\eta}_{\boldsymbol{\theta}}}(s)$  as  $\tilde{\mathbf{p}}_{\boldsymbol{\theta}}(s)$ , then

$$\tilde{\mathbf{p}}_{\boldsymbol{\theta}}(s) = (\tilde{p}_k(s; \boldsymbol{\theta}))_{k=0}^{K-1} = \mathbb{E} \left[ \tilde{\mathbf{G}}(r_0) \mathbf{C}^{-1} \boldsymbol{\Theta}^\top \phi(s_1) \middle| s_0 = s \right] + \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ \mathbf{g}_j(r_0) \middle| s_0 = s \right].$$

Combining this proposition with Proposition 3.1, we know that  $\boldsymbol{\Theta}^*$  is the unique solution to the

following system of linear equations for  $\Theta \in \mathbb{R}^{d \times K}$

$$\begin{aligned}
\Theta &= \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \left( \tilde{p}_\theta(s) - \frac{1}{K+1} \mathbf{1}_K \right)^\top C^\top \right] \\
&= \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \left( \frac{1}{K+1} \mathbf{1}_K - \mathbb{E} \left[ \tilde{G}(r_0) C^{-1} \Theta^\top \phi(s_1) \middle| s_0 = s \right] - \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ g_j(r_0) \middle| s_0 = s \right] \right)^\top C^\top \right] \\
&= \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \phi(s')^\top \Theta \left( C \tilde{G}(r) C^{-1} \right)^\top \right] + \frac{1}{K+1} \Sigma_\phi^{-1} \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \left( \sum_{j=0}^K g_j(r) - \mathbf{1}_K \right)^\top C^\top \right],
\end{aligned}$$

or equivalently,

$$\Sigma_\phi \Theta - \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \phi(s')^\top \Theta \left( C \tilde{G}(r) C^{-1} \right)^\top \right] = \frac{1}{K+1} \mathbb{E}_{s \sim \mu_\pi} \left[ \phi(s) \left( \sum_{j=0}^K g_j(r) - \mathbf{1}_K \right)^\top C^\top \right],$$

which is the desired conclusion. The uniqueness and existence of the solution is guaranteed by the fact that the LHS is an invertible linear transformation of  $\Theta$ , which is justified by Eqn. (28).

#### D.4 Proof of Proposition D.3

*Proof.* Recall the definition of the distributional Bellman operator Eqn. (6) and categorical projection operator Eqn. (8), we have

$$\begin{aligned}
\tilde{p}_k(s; \theta) &= p_k([\mathcal{T}^\pi \eta_\theta](s)) \\
&= \mathbb{E}_{X \sim [\mathcal{T}^\pi \eta_\theta](s)} \left[ \left( 1 - \left| \frac{X - x_k}{\iota_K} \right| \right)_+ \right] \\
&= \mathbb{E} \left[ \mathbb{E}_{G \sim \eta_\theta(s_1)} \left[ \left( 1 - \left| \frac{r_0 + \gamma G - x_k}{\iota_K} \right| \right)_+ \right] \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \sum_{j=0}^K p_j(s_1; \theta) \left( 1 - \left| \frac{r_0 + \gamma x_j - x_k}{\iota_K} \right| \right)_+ \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \sum_{j=0}^K p_j(s_1; \theta) g_{j,k}(r_0) \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ g_{K,k}(r_0) + \sum_{j=0}^{K-1} p_j(s_1; \theta) (g_{j,k}(r_0) - g_{K,k}(r_0)) \middle| s_0 = s \right].
\end{aligned} \tag{20}$$



Hence, let  $\mathbf{W} = \mathbf{\Theta} \mathbf{C}^{-\top}$  and  $\mathbf{w} = \text{vec}(\mathbf{W}) = (\mathbf{C}^{-1} \otimes \mathbf{I}_d) \boldsymbol{\theta}$  (see Appendix D.6 for their meaning), then

$$\begin{aligned}
\tilde{\mathbf{p}}_{\boldsymbol{\theta}}(s) &= (\tilde{p}_k(s; \boldsymbol{\theta}))_{k=0}^{K-1} \\
&= \mathbb{E} \left[ \begin{bmatrix} g_{K,1}(r_0) \\ \vdots \\ g_{K,K-1}(r_0) \end{bmatrix} + \sum_{j=0}^{K-1} p_j(s_1; \boldsymbol{\theta}) \begin{bmatrix} g_{j,1}(r_0) - g_{K,1}(r_0) \\ \vdots \\ g_{j,K-1}(r_0) - g_{K,K-1}(r_0) \end{bmatrix} \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \mathbf{g}_K(r_0) + \sum_{j=0}^{K-1} p_j(s_1; \boldsymbol{\theta}) (\mathbf{g}_j(r_0) - \mathbf{g}_K(r_0)) \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \mathbf{g}_K(r_0) + (\mathbf{G}(r_0) - \mathbf{1}_K^{\top} \otimes \mathbf{g}_K(r_0)) \mathbf{p}_{\boldsymbol{\theta}}(s_1) \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \mathbf{g}_K(r_0) + (\mathbf{G}(r_0) - \mathbf{1}_K^{\top} \otimes \mathbf{g}_K(r_0)) \left[ (\mathbf{I}_K \otimes \boldsymbol{\phi}(s_1))^{\top} \mathbf{w} + \frac{1}{K+1} \mathbf{1}_K \right] \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ (\mathbf{G}(r_0) - \mathbf{1}_K^{\top} \otimes \mathbf{g}_K(r_0)) (\mathbf{I}_K \otimes \boldsymbol{\phi}(s_1))^{\top} \middle| s_0 = s \right] \mathbf{w} \\
&\quad + \mathbb{E} \left[ \mathbf{g}_K(r_0) + \frac{1}{K+1} (\mathbf{G}(r_0) - \mathbf{1}_K^{\top} \otimes \mathbf{g}_K(r_0)) \mathbf{1}_K \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ (\mathbf{G}(r_0) - \mathbf{1}_K^{\top} \otimes \mathbf{g}_K(r_0)) \otimes \boldsymbol{\phi}(s_1)^{\top} \middle| s_0 = s \right] \mathbf{w} + \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ \mathbf{g}_j(r_0) \middle| s_0 = s \right],
\end{aligned}$$

or equivalently,

$$\begin{aligned}
\tilde{\mathbf{p}}_{\boldsymbol{\theta}}(s) &= \mathbb{E} \left[ \tilde{\mathbf{G}}(r_0) \mathbf{W}^{\top} \boldsymbol{\phi}(s_1) \middle| s_0 = s \right] + \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ \mathbf{g}_j(r_0) \middle| s_0 = s \right] \\
&= \mathbb{E} \left[ \tilde{\mathbf{G}}(r_0) \mathbf{C}^{-1} \mathbf{\Theta}^{\top} \boldsymbol{\phi}(s_1) \middle| s_0 = s \right] + \frac{1}{K+1} \sum_{j=0}^K \mathbb{E} \left[ \mathbf{g}_j(r_0) \middle| s_0 = s \right].
\end{aligned}$$

□

## D.5 Proof of Proposition 3.2

*Proof.* By the basic inequality (Lemma G.1), we only need to show

$$\begin{aligned}\ell_{2,\mu_\pi}^2(\boldsymbol{\eta}^\pi, \boldsymbol{\eta}_{\theta^\star}) &\leq \frac{\ell_{2,\mu_\pi}^2(\boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_{\phi,K}^\pi \boldsymbol{\eta}^\pi)}{1-\gamma} \\ &= \frac{\ell_{2,\mu_\pi}^2(\boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_K \boldsymbol{\eta}^\pi) + \ell_{2,\mu_\pi}^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_{\phi,K}^\pi \boldsymbol{\eta}^\pi)}{1-\gamma} \\ &\leq \frac{1}{K(1-\gamma)^2} + \frac{\ell_{2,\mu_\pi}^2(\boldsymbol{\Pi}_K \boldsymbol{\eta}^\pi, \boldsymbol{\Pi}_{\phi,K}^\pi \boldsymbol{\eta}^\pi)}{1-\gamma},\end{aligned}$$

where we used Bellemare et al. [4, Proposition 9.18 and Eqn. (5.28)].  $\square$

## D.6 Linear-Categorical Parametrization with Probability Mass Function Representation

We introduce new notations for linear-categorical parametrization with probability mass function (PMF) representation. Let  $\mathbf{W} := \boldsymbol{\Theta} \mathbf{C}^{-\top} = (\boldsymbol{\theta}(0), \boldsymbol{\theta}(1) - \boldsymbol{\theta}(0), \dots, \boldsymbol{\theta}(K-1) - \boldsymbol{\theta}(K-2)) \in \mathbb{R}^{d \times K}$  and the vectorization of  $\mathbf{W}$ ,  $\mathbf{w} := \text{vec}(\mathbf{W}) = (\mathbf{C}^{-1} \otimes \mathbf{I}_d) \boldsymbol{\theta} \in \mathbb{R}^{dK}$ . We abbreviate  $\mathbf{p}_{\boldsymbol{\eta}_\theta}$  as  $\mathbf{p}_\mathbf{w}$  in this section. Then by Lemma A.2, for any  $s \in \mathcal{S}$ , it holds that

$$\mathbf{p}_\mathbf{w}(s) = \mathbf{W}^\top \boldsymbol{\phi}(s) + (K+1)^{-1} \mathbf{1}_K. \quad (21)$$

PMF and CDF representations are equivalent because  $\mathbf{C}$  is invertible.

For any  $\boldsymbol{\eta}_{\mathbf{w}_1}, \boldsymbol{\eta}_{\mathbf{w}_2} \in \mathcal{P}_{\phi,K}^{\text{sign}}$ , by Proposition C.1,

$$\begin{aligned}\ell_{2,\mu_\pi}^2(\boldsymbol{\eta}_{\mathbf{w}_1}, \boldsymbol{\eta}_{\mathbf{w}_2}) &= \mathbb{E}_{s \sim \mu_\pi} [\ell_2^2(\boldsymbol{\eta}_{\mathbf{w}_1}(s), \boldsymbol{\eta}_{\mathbf{w}_2}(s))] \\ &= \iota_K \mathbb{E}_{s \sim \mu_\pi} [\|\mathbf{C}(\mathbf{p}_{\mathbf{w}_1}(s) - \mathbf{p}_{\mathbf{w}_2}(s))\|^2] \\ &= \iota_K \text{tr} \left( \boldsymbol{\Sigma}_\phi^{\frac{1}{2}} (\mathbf{W}_1 - \mathbf{W}_2) \mathbf{C}^\top \mathbf{C} (\mathbf{W}_1 - \mathbf{W}_2)^\top \boldsymbol{\Sigma}_\phi^{\frac{1}{2}} \right) \\ &= \iota_K \|\mathbf{w}_1 - \mathbf{w}_2\|_{(\mathbf{C}^\top \mathbf{C}) \otimes \boldsymbol{\Sigma}_\phi}^2,\end{aligned} \quad (22)$$

hence the affine space  $(\mathcal{P}_{\phi,K}^{\text{sign}}, \ell_{2,\mu_\pi})$  is isometric with the Euclidean space  $(\mathbb{R}^{dK}, \sqrt{\iota_K} \|\cdot\|_{(\mathbf{C}^\top \mathbf{C}) \otimes \boldsymbol{\Sigma}_\phi})$  if we consider the PMF representation.

Following the proof of Lemma D.1 in Appendix D.2, we can also derive the gradient when we

use the PMF parametrization:

$$\begin{aligned}
\nabla_{\mathbf{w}} \ell_2^2(\eta_{\mathbf{w}}(s), \eta(s)) &= \nabla_{\mathbf{w}} \ell_2^2(\eta_{\mathbf{w}}(s), \mathbf{\Pi}_K \eta(s)) \\
&= \iota_K \nabla_{\mathbf{w}} \|\mathbf{C}(\mathbf{p}_{\mathbf{w}}(s) - \mathbf{p}_{\eta}(s))\|^2 \\
&= 2\iota_K (\mathbf{I}_K \otimes \phi(s)) \mathbf{C}^\top \mathbf{C} (\mathbf{p}_{\mathbf{w}}(s) - \mathbf{p}_{\eta}(s)) \\
&= 2\iota_K (\mathbf{I}_K \otimes \phi(s)) \mathbf{C}^\top \mathbf{C} \left( (\mathbf{I}_K \otimes \phi(s))^\top \mathbf{w} + \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\eta}(s) \right) \\
&= 2\iota_K \left\{ [(\mathbf{C}^\top \mathbf{C}) \otimes (\phi(s)\phi(s)^\top)] \mathbf{w} + \left[ \left( \mathbf{C}^\top \mathbf{C} \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\eta}(s) \right) \right) \otimes \phi(s) \right] \right\}, \\
\\
\nabla_{\mathbf{W}} \ell_2^2(\eta_{\mathbf{w}}(s), \eta(s)) &= 2\iota_K \phi(s) (\mathbf{p}_{\mathbf{w}}(s) - \mathbf{p}_{\eta}(s))^\top \mathbf{C}^\top \mathbf{C} \\
&= 2\iota_K \phi(s) \left[ \phi(s)^\top \mathbf{W} + \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\eta}(s) \right)^\top \right] \mathbf{C}^\top \mathbf{C}. \tag{23}
\end{aligned}$$

## D.7 Stochastic Semi-Gradient Descent with Linear Function Approximation

We denote by  $\mathcal{T}_t^\pi$  the corresponding empirical distributional Bellman operator at the  $t$ -th iteration, which satisfies

$$[\mathcal{T}_t^\pi \eta](s_t) = (b_{r_t, \gamma})_\#(\eta(s_{t+1})), \quad \forall \eta \in \mathcal{P}^{\mathcal{S}}. \tag{24}$$

### D.7.1 Probability Mass Function Representation

Consider the stochastic semi-gradient descent (SSGD) with the probability mass function (PMF) representation

$$\mathbf{W}_t \leftarrow \mathbf{W}_{t-1} - \alpha \nabla_{\mathbf{W}} \ell_2^2(\eta_{\mathbf{w}_{t-1}}(s_t), [\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}](s_t)),$$

where  $\nabla_{\mathbf{W}}$  stands for taking gradient w.r.t.  $\mathbf{W}_{t-1} \in \mathbb{R}^{d \times K}$  in the first term  $\eta_{\mathbf{w}_{t-1}}(s_t)$  (the second term is regarding as a constant, that's why we call it a semi-gradient). We can check that  $\nabla_{\mathbf{W}} \ell_2^2(\eta_{\mathbf{w}_{t-1}}(s_t), [\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}](s_t))$  is an unbiased estimate of  $\nabla_{\mathbf{W}} \ell_{2, \mu_\pi}^2(\eta_{\mathbf{w}_{t-1}}, \mathcal{T}^\pi \eta_{\mathbf{w}_{t-1}})$ .

Now, let us compute the gradient term. By Eqn. (23), we have

$$\begin{aligned}
\nabla_{\mathbf{W}} \ell_2^2(\eta_{\mathbf{w}_{t-1}}(s_t), [\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}](s_t)) &= 2\iota_K \phi(s_t) \left( \mathbf{p}_{\mathbf{w}_{t-1}}(s_t) - \mathbf{p}_{\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}}(s_t) \right)^\top \mathbf{C}^\top \mathbf{C} \\
&= 2\iota_K \phi(s_t) \left[ \phi(s_t)^\top \mathbf{W}_{t-1} + \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}}(s_t) \right)^\top \right] \mathbf{C}^\top \mathbf{C}.
\end{aligned}$$

where  $\mathbf{p}_{\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}}(s_t) = \mathbf{p}_{\mathbf{\Pi}_K \mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}(s_t)} = (p_k([\mathcal{T}_t^\pi \eta_{\mathbf{w}_{t-1}}](s_t)))_{k=0}^{K-1} \in \mathbb{R}^K$ . Now, we turn to compute

$\mathbf{p}\mathcal{T}_t^\pi \boldsymbol{\eta}_{\mathbf{w}_{t-1}}(s_t)$ . According to Eqn. (8),

$$\begin{aligned}
p_k([\mathcal{T}_t^\pi \boldsymbol{\eta}_{\mathbf{w}_{t-1}}](s_t)) &= \mathbb{E}_{X \sim [\mathcal{T}_t^\pi \boldsymbol{\eta}_{\mathbf{w}_{t-1}}](s_t)} \left[ \left( 1 - \left| \frac{X - x_k}{\iota_K} \right| \right)_+ \right] \\
&= \mathbb{E}_{G \sim \boldsymbol{\eta}_{\mathbf{w}_{t-1}}(s_{t+1})} \left[ \left( 1 - \left| \frac{r_t + \gamma G - x_k}{\iota_K} \right| \right)_+ \right] \\
&= \sum_{j=0}^K p_j(s_{t+1}; \mathbf{w}_{t-1}) g_{j,k}(r_t) \\
&= g_{K,k}(r_t) + \sum_{j=0}^{K-1} p_j(s_{t+1}; \mathbf{w}_{t-1}) (g_{j,k}(r_t) - g_{K,k}(r_t)),
\end{aligned}$$

which has the same form as Eqn. (20). Following the proof of Proposition D.3 in Appendix D.4, one can show that

$$\mathbf{p}\mathcal{T}_t^\pi \boldsymbol{\eta}_{\mathbf{w}_{t-1}}(s_t) = \tilde{\mathbf{G}}(r_t) \mathbf{W}_{t-1}^\top \boldsymbol{\phi}(s_{t+1}) + \frac{1}{K+1} \sum_{j=0}^K \mathbf{g}_j(r_t). \quad (25)$$

Hence, the update scheme is

$$\begin{aligned}
\mathbf{W}_t &\leftarrow \mathbf{W}_{t-1} - 2\iota_K \alpha \boldsymbol{\phi}(s_t) \left( \mathbf{p}_{\mathbf{w}_{t-1}}(s_t) - \mathbf{p}\mathcal{T}_t^\pi \boldsymbol{\eta}_{\mathbf{w}_{t-1}}(s_t) \right)^\top \mathbf{C}^\top \mathbf{C} \\
&= \mathbf{W}_{t-1} - 2\iota_K \alpha \boldsymbol{\phi}(s_t) \left[ \boldsymbol{\phi}(s_t)^\top \mathbf{W}_{t-1} - \boldsymbol{\phi}(s_{t+1})^\top \mathbf{W}_{t-1} \tilde{\mathbf{G}}^\top(r_t) - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^\top \right] \mathbf{C}^\top \mathbf{C}.
\end{aligned} \quad (26)$$

Note that our Linear-CTD (Eqn. (13)) is equivalent to

$$\mathbf{W}_t \leftarrow \mathbf{W}_{t-1} - \alpha \boldsymbol{\phi}(s_t) \left[ \boldsymbol{\phi}(s_t)^\top \mathbf{W}_{t-1} - \boldsymbol{\phi}(s_{t+1})^\top \mathbf{W}_{t-1} \tilde{\mathbf{G}}^\top(r_t) - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^\top \right], \quad (27)$$

in the PMF representation. Compared to Eqn. (27), the SSGD (Eqn. (26)) has an additional  $\mathbf{C}^\top \mathbf{C}$ , and the step size is  $2\iota_K \alpha$ .

### D.7.2 Cumulative Distribution Function Representation

Consider the SSGD with the CDF representation

$$\boldsymbol{\Theta}_t \leftarrow \boldsymbol{\Theta}_{t-1} - \alpha \nabla_{\boldsymbol{\Theta}} \ell_2^2(\boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}(s_t), [\mathcal{T}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}](s_t)),$$

where  $\nabla_{\boldsymbol{\Theta}}$  stands for taking gradient w.r.t.  $\boldsymbol{\Theta}_{t-1} = \boldsymbol{\theta}_{t-1} \mathbf{C}^\top \in \mathbb{R}^{d \times K}$  in the first term  $\boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}(s_t)$  (the second term is regarding as a constant). One can check that  $\nabla_{\boldsymbol{\Theta}} \ell_2^2(\boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}(s_t), [\mathcal{T}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}](s_t))$  is an

unbiased estimate of  $\nabla_{\Theta} \ell_{2, \mu_{\pi}}^2(\eta_{\theta_{t-1}}, \mathcal{T}^{\pi} \eta_{\theta_{t-1}})$ .

Now, let us compute the gradient term. By Eqn. (19) and Eqn. (25), we have

$$\begin{aligned} \nabla_{\Theta} \ell_2^2(\eta_{\theta_{t-1}}(s_t), [\mathcal{T}_t^{\pi} \eta_{\theta_{t-1}}](s_t)) &= 2\iota_K \phi(s_t) \left( \mathbf{F}_{\theta_{t-1}}(s_t) - \mathbf{F}_{\mathcal{T}_t^{\pi} \eta_{\theta_{t-1}}}(s_t) \right)^{\top} \\ &= 2\iota_K \phi(s_t) \left[ \phi(s_t)^{\top} \Theta_{t-1} + \left( \frac{1}{K+1} \mathbf{1}_K - \mathbf{p}_{\mathcal{T}_t^{\pi} \eta_{\theta_{t-1}}}(s) \right)^{\top} \mathbf{C}^{\top} \right] \\ &= 2\iota_K \phi(s_t) \left[ \phi(s_t)^{\top} \Theta_{t-1} - \phi(s_{t+1})^{\top} \Theta_{t-1} \mathbf{C}^{-\top} \tilde{\mathbf{G}}^{\top}(r_t) \mathbf{C}^{\top} - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^{\top} \mathbf{C}^{\top} \right]. \end{aligned}$$

Hence, the update scheme is

$$\begin{aligned} \Theta_t &\leftarrow \Theta_{t-1} - 2\iota_K \alpha \phi(s_t) \left( \mathbf{F}_{\theta_{t-1}}(s_t) - \mathbf{F}_{\mathcal{T}_t^{\pi} \eta_{\theta_{t-1}}}(s_t) \right)^{\top} \\ &= \Theta_{t-1} - 2\iota_K \alpha \phi(s_t) \left[ \phi(s_t)^{\top} \Theta_{t-1} - \phi(s_{t+1})^{\top} \Theta_{t-1} \mathbf{C}^{-\top} \tilde{\mathbf{G}}^{\top}(r_t) \mathbf{C}^{\top} - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^{\top} \mathbf{C}^{\top} \right], \end{aligned}$$

which has the same form as Linear-CTD (Eqn. (13)) with the step size  $2\alpha\iota_K$ .

## D.8 Linear-CTD is mean-preserving

We will show that our Linear-CTD is mean-preserving, which was first discovered by Lyle et al. [32, Proposition 8]. In this section, we assume the first coordinate of the feature is a constant, i.e.,  $\phi_1(s) = 1/\sqrt{d}$  for any  $s \in \mathcal{S}$ . As stated before, we will always assume this to ensure  $\mathcal{P}_{\phi, K}$  can be uniquely defined.

**Proposition D.4.** *Let  $\mathbf{V}_{\theta} := (V_{\theta}(s))_{s \in \mathcal{S}}$  be the value function corresponding to  $\theta$ , i.e.,  $V_{\theta}(s) = \mathbb{E}_{G \sim \eta_{\theta}(s)}[G]$ , then for any initialization of the Linear-TD parameter  $\psi_0$ , there exists a (not unique) corresponding Linear-CTD parameter  $\theta_0$  such that  $\mathbf{V}_{\theta_0} = \mathbf{V}_{\psi_0}$ , furthermore, for any  $t \geq 1$  and even number  $T \geq 2$ , it holds that*

$$\mathbf{V}_{\theta_t} = \mathbf{V}_{\psi_t}, \quad \mathbf{V}_{\theta_T} = \mathbf{V}_{\tilde{\psi}_T}.$$

*Proof of Proposition D.4.* Recall that the updating scheme of Linear-TD is given by

$$\psi_t \leftarrow \psi_{t-1} - \alpha \phi(s_t) (V_{\psi_{t-1}}(s_t) - r_t - \gamma V_{\psi_{t-1}}(s_{t+1}))$$

And the updating scheme of Linear-CTD is given by

$$\boldsymbol{\Theta}_t \leftarrow \boldsymbol{\Theta}_{t-1} - \alpha \phi(s_t) \left( \mathbf{F}_{\boldsymbol{\Theta}_{t-1}}(s_t) - \mathbf{F}_{\mathcal{T}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}}(s_t) \right)^\top.$$

Let  $\mathbf{V}_\boldsymbol{\theta} := (V_\boldsymbol{\theta}(s))_{s \in \mathcal{S}}$  be the value function corresponding to  $\boldsymbol{\theta}$ , we have

$$\begin{aligned} V_\boldsymbol{\theta}(s) &= \iota_K \sum_{k=0}^{K-1} (1 - F_k(s; \boldsymbol{\theta})) \\ &= \iota_K (\mathbf{1}_K - \mathbf{F}_\boldsymbol{\theta}(s))^\top \mathbf{1}_K \\ &= \frac{1}{2(1-\gamma)} - \iota_K \phi(s)^\top \boldsymbol{\Theta} \mathbf{1}_K. \end{aligned}$$

Hence, if we take  $\psi_{0,1} = \frac{\sqrt{d}}{2(1-\gamma)}$ ,  $\psi_{0,i} = 0$  for any  $i \in \{2, \dots, d\}$ , and  $\boldsymbol{\theta}_0 = \mathbf{0}_{dK}$ , it holds that

$$\mathbf{V}_{\psi_0} = \mathbf{V}_{\boldsymbol{\theta}_0} = \frac{1}{2(1-\gamma)} \mathbf{1}_\mathcal{S}.$$

We can also show that for any  $\psi_0 \in \mathbb{R}^d$ , there exists  $\boldsymbol{\theta}_0 \in \mathbb{R}^{d \times K}$  such that  $\mathbf{V}_{\psi_0} = \mathbf{V}_{\boldsymbol{\theta}_0}$ . That is we need to find  $\boldsymbol{\theta}_0$  such that for any  $s \in \mathcal{S}$ ,

$$\iota_K \sum_{k=0}^{K-1} \phi(s)^\top \boldsymbol{\theta}_0(k) = \frac{1}{2(1-\gamma)} - \phi(s)^\top \psi_0.$$

We can take  $\boldsymbol{\theta}_0$  satisfying the following equations to make the above equation hold

$$\iota_K \sum_{k=0}^{K-1} \theta_0(k, 1) = \frac{\sqrt{d}}{2(1-\gamma)} - \psi_{0,1},$$

and

$$\iota_K \sum_{k=0}^{K-1} \boldsymbol{\theta}_0(k)_{-1} = -\psi_{0,-1}.$$

Furthermore, for any  $t \geq 1$ , we have for any  $s \in \mathcal{S}$ ,

$$\begin{aligned} V_{\boldsymbol{\theta}_t}(s) &= \frac{1}{2(1-\gamma)} - \iota_K \phi(s)^\top \boldsymbol{\Theta}_t \mathbf{1}_K \\ &= \frac{1}{2(1-\gamma)} - \iota_K \phi(s)^\top \left( \boldsymbol{\Theta}_{t-1} - \alpha \phi(s_t) \left( \mathbf{F}_{\boldsymbol{\Theta}_{t-1}}(s_t) - \mathbf{F}_{\mathcal{T}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}}(s_t) \right)^\top \right) \mathbf{1}_K \\ &= V_{\boldsymbol{\Theta}_{t-1}}(s) + \alpha \iota_K (\phi(s)^\top \phi(s_t)) \left( \mathbf{F}_{\boldsymbol{\Theta}_{t-1}}(s_t) - \mathbf{F}_{\mathcal{T}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\Theta}_{t-1}}}(s_t) \right)^\top \mathbf{1}_K \end{aligned}$$

$$\begin{aligned}
V_{\psi_t}(s) &= \phi(s)^\top \psi_t \\
&= \phi(s)^\top (\psi_{t-1} - \alpha \phi(s_t) (V_{\psi_{t-1}}(s_t) - r_t - \gamma V_{\psi_{t-1}}(s_{t+1}))) \\
&= V_{\psi_{t-1}}(s) - \alpha (\phi(s)^\top \phi(s_t)) (V_{\psi_{t-1}}(s_t) - r_t - \gamma V_{\psi_{t-1}}(s_{t+1})).
\end{aligned}$$

We need to check that, if  $\mathbf{V}_{\boldsymbol{\theta}_{t-1}} = \mathbf{V}_{\psi_{t-1}}$ , it holds that

$$\iota_K \left( \mathbf{F} \boldsymbol{\tau}_t^\pi \boldsymbol{\eta}_{\boldsymbol{\theta}_{t-1}}(s_t) - \mathbf{F}_{\boldsymbol{\theta}_{t-1}}(s_t) \right)^\top \mathbf{1}_K = V_{\psi_{t-1}}(s_t) - r_t - \gamma V_{\psi_{t-1}}(s_{t+1}),$$

which is the direct corollary of the following fact

$$\text{LHS} = V_{\psi_t}(s_t) - \mathbb{E}_{X \sim \Pi_K(b_{r,\gamma})_{\#} \boldsymbol{\eta}_{\boldsymbol{\theta}_{t-1}}(s_{t+1})}[X] = V_{\psi_t}(s_t) - r_t - \gamma V_{\boldsymbol{\theta}_{t-1}}(s_{t+1}),$$

by Lemma D.2. And we can obtain  $\mathbf{V}_{\bar{\boldsymbol{\theta}}_T} = \mathbf{V}_{\bar{\psi}_T}$  by using the facts that  $\mathcal{P}_{\phi,K}^{\text{sign}}$  is an affine space,  $\mathcal{V}_\phi$  is a linear space and taking expectation is a linear operator.  $\square$

**Lemma D.2.** *For any  $\nu \in \mathcal{P}^{\text{sign}}$ , it holds that*

$$\mathbb{E}_{X \sim \nu}[X] = \mathbb{E}_{X \sim \Pi_K \nu}[X].$$

*Proof.* By Eqn. (8),  $\Pi_K \nu \in \mathcal{P}_K^{\text{sign}}$  is uniquely identified with a vector  $\mathbf{p}_\nu = (p_k(\nu))_{k=0}^{K-1} \in \mathbb{R}^K$ , where

$$p_k(\nu) = \int_{[0, (1-\gamma)^{-1}]} (1 - |(x - x_k)/\iota_K|)_+ \nu(dx).$$

Hence, for any  $x \in [0, (1-\gamma)^{-1}]$ , we define  $x_{\text{lb}} := \max\{y \in \{x_0, \dots, x_K\} : x \leq y\}$ , then

$$\begin{aligned}
\sum_{k=0}^K x_k (1 - |(x - x_k)/\iota_K|)_+ &= x_{\text{lb}} \left( 1 - \frac{x - x_{\text{lb}}}{\iota_K} \right) + (x_{\text{lb}} + \iota_K) \left( 1 - \frac{x_{\text{lb}} + \iota_K - x}{\iota_K} \right) \\
&= x_{\text{lb}} + \iota_K \left( 1 - \frac{x_{\text{lb}} + \iota_K - x}{\iota_K} \right) \\
&= x,
\end{aligned}$$

therefore,

$$\begin{aligned}
\mathbb{E}_{X \sim \Pi_K \nu}[X] &= \sum_{k=0}^K x_k \int_{[0, (1-\gamma)^{-1}]} (1 - |(x - x_k)/\iota_K|)_+ \nu(dx) \\
&= \int_{[0, (1-\gamma)^{-1}]} \sum_{k=0}^K x_k (1 - |(x - x_k)/\iota_K|)_+ \nu(dx) \\
&= \int_{[0, (1-\gamma)^{-1}]} x \nu(dx) \\
&= \mathbb{E}_{X \sim \nu}[X].
\end{aligned}$$

□

## E Omitted Results and Proofs in Section 4

### E.1 Proof of Lemma 5.1

*Proof.* By Lemma G.1 and Eqn. (22), we have

$$\begin{aligned}
(\mathcal{L}(\boldsymbol{\theta}))^2 &= W_{1, \mu_\pi}^2(\boldsymbol{\eta}_\theta, \boldsymbol{\eta}_{\theta^*}) \\
&\leq \frac{1}{1-\gamma} \ell_{2, \mu_\pi}^2(\boldsymbol{\eta}_\theta, \boldsymbol{\eta}_{\theta^*}) \\
&= \frac{\iota_K}{1-\gamma} \text{tr} \left( (\boldsymbol{\Theta} - \boldsymbol{\Theta}^*)^\top \boldsymbol{\Sigma}_\phi (\boldsymbol{\Theta} - \boldsymbol{\Theta}^*) \right) \\
&= \frac{1}{K(1-\gamma)^2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_{\boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi}^2,
\end{aligned}$$

We only need to show that

$$\boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi \leq \frac{1}{(1-\sqrt{\gamma})^2} \bar{\boldsymbol{A}}^\top \left( \boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-1} \right) \bar{\boldsymbol{A}} \left( \leq \frac{4}{(1-\gamma)^2 \lambda_{\min}} \bar{\boldsymbol{A}}^\top \bar{\boldsymbol{A}} \right), \quad (28)$$

or equivalently,

$$\left( \boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \bar{\boldsymbol{A}}^\top \left( \boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-1} \right) \bar{\boldsymbol{A}} \left( \boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \geq (1-\sqrt{\gamma})^2 \boldsymbol{I}_{dK}.$$

Recall

$$\bar{\boldsymbol{A}} = (\boldsymbol{I}_K \otimes \boldsymbol{\Sigma}_\phi) - \mathbb{E}_{s, r, s'} \left[ \left( C \tilde{G}(r) C^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right],$$



then for any  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$  with  $\|\boldsymbol{\theta}\| = 1$ ,

$$\begin{aligned}
& \boldsymbol{\theta}^\top \left( \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \bar{\mathbf{A}}^\top \left( \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-1} \right) \bar{\mathbf{A}} \left( \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \boldsymbol{\theta} \\
&= \left\| \left( \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \bar{\mathbf{A}} \left( \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \boldsymbol{\theta} \right\|^2 \\
&= \left\| \boldsymbol{\theta} - \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \boldsymbol{\theta} \right\|^2 \\
&\geq \left( 1 - \left\| \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \boldsymbol{\theta} \right\| \right)^2.
\end{aligned}$$

It suffices to show that

$$\left\| \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \right\| \leq \sqrt{\gamma}. \quad (29)$$

For brevity, we abbreviate  $\mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1}$  as  $\mathbf{Y}(r) = (y_{ij}(r))_{i,j=1}^K \in \mathbb{R}^{K \times K}$ . Thus, it suffices to show that

$$\left\| \mathbb{E}_{s,r,s'} \left[ \mathbf{Y}(r) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \right\| \leq \sqrt{\gamma}.$$

For any vectors  $\mathbf{w} = (\mathbf{w}(0)^\top, \dots, \mathbf{w}(K-1)^\top)$  and  $\mathbf{v} = (\mathbf{v}(0)^\top, \dots, \mathbf{v}(K-1)^\top)$  in  $\mathbb{R}^{dK}$ , we define the corresponding matrices  $\mathbf{W} = (\mathbf{w}(0), \dots, \mathbf{w}(K-1))$  and  $\mathbf{V} = (\mathbf{v}(0), \dots, \mathbf{v}(K-1))$  in  $\mathbb{R}^{d \times K}$ , then  $\|\mathbf{w}\| = \|\mathbf{W}\|_F = \sqrt{\text{tr}(\mathbf{W}^\top \mathbf{W})}$  and  $\|\mathbf{v}\| = \|\mathbf{V}\|_F = \sqrt{\text{tr}(\mathbf{V}^\top \mathbf{V})}$ . With these notations, we have

$$\begin{aligned}
& \left\| \mathbb{E}_{s,r,s'} \left[ \mathbf{Y}(r) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \right\| \\
&= \sup_{\|\mathbf{w}\|=\|\mathbf{v}\|=1} \mathbf{w}^\top \mathbb{E}_{s,r,s'} \left[ \mathbf{Y}(r) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \mathbf{v} \\
&= \sup_{\|\mathbf{w}\|=\|\mathbf{v}\|=1} \mathbb{E}_{s,r,s'} \left[ \sum_{i,j=1}^K y_{ij}(r) \mathbf{w}(i)^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \mathbf{v}(j) \right],
\end{aligned}$$

it is easy to check that

$$\mathbf{w}(i)^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \mathbf{v}(j) = \left( \mathbf{W}^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \mathbf{V} \right)_{ij},$$

hence

$$\begin{aligned}
& \left\| \mathbb{E}_{s,r,s'} \left[ \mathbf{Y}(r) \otimes \left( \Sigma_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \Sigma_\phi^{-\frac{1}{2}} \right) \right] \right\| \\
&= \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \mathbb{E}_{s,r,s'} \left[ \sum_{i,j=1}^K y_{ij}(r) \left( \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \Sigma_\phi^{-\frac{1}{2}} \mathbf{V} \right)_{ij} \right] \\
&= \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \mathbb{E}_{s,r,s'} \left[ \text{tr} \left( \mathbf{Y}(r) \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \phi(s)^\top \Sigma_\phi^{-\frac{1}{2}} \mathbf{W} \right) \right] \\
&= \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \mathbb{E}_{s,r,s'} \left[ \phi(s)^\top \Sigma_\phi^{-\frac{1}{2}} \mathbf{W} \mathbf{Y}(r) \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \right] \\
&\leq \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \mathbb{E}_{s,r,s'} \left[ \left\| \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s) \right\| \left\| \mathbf{Y}(r) \right\| \left\| \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \right\| \right] \\
&\leq \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \mathbb{E}_{s,s'} \left[ \left\| \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s) \right\| \left\| \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \right\| \right] \\
&\leq \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \sqrt{\mathbb{E}_s \left[ \left\| \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s) \right\|^2 \right] \mathbb{E}_{s'} \left[ \left\| \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \right\|^2 \right]} \\
&= \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \sqrt{\mathbb{E}_s \left[ \phi(s)^\top \Sigma_\phi^{-\frac{1}{2}} \mathbf{W} \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s) \right] \mathbb{E}_{s'} \left[ \phi(s')^\top \Sigma_\phi^{-\frac{1}{2}} \mathbf{V} \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \phi(s') \right]} \\
&= \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \sqrt{\text{tr} \left( \mathbf{W} \mathbf{W}^\top \Sigma_\phi^{-\frac{1}{2}} \mathbb{E}_s [\phi(s) \phi(s)^\top] \Sigma_\phi^{-\frac{1}{2}} \right) \text{tr} \left( \mathbf{V} \mathbf{V}^\top \Sigma_\phi^{-\frac{1}{2}} \mathbb{E}_{s'} [\phi(s') \phi(s')^\top] \Sigma_\phi^{-\frac{1}{2}} \right)} \\
&= \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \sqrt{\text{tr}(\mathbf{W} \mathbf{W}^\top) \text{tr}(\mathbf{V} \mathbf{V}^\top)} \\
&= \sqrt{\gamma} \sup_{\|\mathbf{W}\|_F = \|\mathbf{V}\|_F = 1} \|\mathbf{W}\|_F \|\mathbf{V}\|_F \\
&= \sqrt{\gamma},
\end{aligned}$$

where we used  $\|\mathbf{Y}(r)\| \leq \sqrt{\gamma}$  for any  $r \in [0, 1]$  by Lemma G.3, and Cauchy-Schwarz inequality.

To summarize, we have shown the desired result

$$\left( \mathbf{I}_K \otimes \Sigma_\phi^{-\frac{1}{2}} \right) \bar{\mathbf{A}}^\top \left( \mathbf{I}_K \otimes \Sigma_\phi^{-1} \right) \bar{\mathbf{A}} \left( \mathbf{I}_K \otimes \Sigma_\phi^{-\frac{1}{2}} \right) \geq (1 - \sqrt{\gamma})^2 \mathbf{I}_{dK}.$$

□

## E.2 Proof of Lemma 5.2

*Proof.* For simplicity, we omit  $t$  in the random variables, for example, we use  $\mathbf{A}$  to denote a random matrix with the same distribution as  $\mathbf{A}_t$ . In addition, we omit the subscripts in the expectation, the

involving random variables are  $s \sim \mu_\pi, a \sim \pi(\cdot | s), (r, s') \sim \mathcal{P}(\cdot, \cdot | s, a)$ .

**Bounding  $C_A$ .** By Lemma A.3,

$$\begin{aligned}\|\mathbf{A}\| &\leq \|\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)\| + \|(C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top)\| \\ &= \|\phi(s)\|^2 + \|\phi(s)\| \|\phi(s')\| \|C\tilde{G}(r)C^{-1}\| \\ &\leq 1 + \sqrt{\gamma},\end{aligned}$$

where we used Lemma G.3. Hence,  $C_A \leq 2(1 + \sqrt{\gamma})$ .

**Bounding  $C_e$ .** By Eqn. (32),

$$\begin{aligned}\|\mathbf{A}\boldsymbol{\theta}^\star\|^2 &= (\boldsymbol{\theta}^\star)^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\theta}^\star \\ &\leq 2 \left( \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)}^2 + \gamma \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes (\phi(s')\phi(s')^\top)}^2 \right) \\ &\leq 2(1 + \gamma) \sup_{s \in \mathcal{S}} \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)}^2 \\ &\leq 2(1 + \gamma) \sup_{s \in \mathcal{S}} \|\phi(s)\|^2 \|\boldsymbol{\theta}^\star\|^2 \\ &\leq 2(1 + \gamma) \|\boldsymbol{\theta}^\star\|^2.\end{aligned}$$

Hence

$$\|\mathbf{A}\boldsymbol{\theta}^\star\| \leq \sqrt{2(1 + \gamma)} \|\boldsymbol{\theta}^\star\|.$$

As for  $\|\mathbf{b}\|$ ,

$$\begin{aligned}\|\mathbf{b}\| &= \frac{1}{K+1} \left\| \left[ \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) \right] \otimes \phi(s) \right\| \\ &\leq \frac{1}{K+1} \left\| \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) \right\| \|\phi(s)\| \\ &\leq \frac{1}{K+1} \left\| \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) \right\|.\end{aligned}\tag{30}$$

By Proposition C.3 with  $\boldsymbol{\eta} \in \mathcal{P}_K^{\text{sign}}$  satisfying  $\eta(\tilde{s}) = \nu$  for all  $\tilde{s} \in \mathcal{S}$ , where  $\nu = (K+1)^{-1} \sum_{k=0}^K \delta_{x_k}$

is the discrete uniform distribution, we can derive that, for any  $r \in [0, 1]$  and  $s' \in \mathcal{S}$ , it holds that

$$\begin{aligned} \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) &= \left( \mathbf{p}_{(b_{r,\gamma})_{\#}\eta(s')} - \frac{1}{K+1} \mathbf{1}_K \right) - \tilde{\mathbf{G}}(r) \left( \mathbf{p}_{\eta}(s') - \frac{1}{K+1} \mathbf{1}_K \right) \\ &= \mathbf{p}_{(b_{r,\gamma})_{\#}\nu} - \frac{1}{K+1} \mathbf{1}_K, \end{aligned}$$

Hence,

$$\begin{aligned} \|\mathbf{b}\| &\leq \frac{1}{K+1} \left\| \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) \right\| \\ &= \left\| \mathbf{C} \left( \mathbf{p}_{(b_{r,\gamma})_{\#}\nu} - \frac{1}{K+1} \mathbf{1}_K \right) \right\| \\ &= \frac{1}{\sqrt{\iota_K}} \ell_2(\mathbf{\Pi}_K(b_{r,\gamma})_{\#}(\nu), \nu) \\ &\leq \sqrt{K(1-\gamma)} \ell_2((b_{r,\gamma})_{\#}(\nu), \nu) \\ &\leq 3\sqrt{K}(1-\gamma), \end{aligned}$$

where we used the orthogonal decomposition (Proposition C.2) and an upper bound for  $\ell_2((b_{r,\gamma})_{\#}(\nu), \nu)$  (Lemma G.4).

In summary,

$$\begin{aligned} \|\mathbf{e}\| &= \|\mathbf{A}\boldsymbol{\theta}^* - \mathbf{b}\| \\ &\leq \|\mathbf{A}\boldsymbol{\theta}^*\| + \|\mathbf{b}\| \\ &\leq \sqrt{2(1+\gamma)} \|\boldsymbol{\theta}^*\| + 3\sqrt{K}(1-\gamma). \end{aligned}$$

Hence,  $C_e \leq \sqrt{2(1+\gamma)} \|\boldsymbol{\theta}^*\| + 3\sqrt{K}(1-\gamma)$ .

**Bounding  $\text{tr}(\boldsymbol{\Sigma}_e)$ .**

$$\text{tr}(\boldsymbol{\Sigma}_e) = \mathbb{E} \left[ \|\mathbf{A}\boldsymbol{\theta}^* - \mathbf{b}\|^2 \right] \leq 2(\boldsymbol{\theta}^*)^\top \mathbb{E}[\mathbf{A}^\top \mathbf{A}] \boldsymbol{\theta}^* + 2\mathbb{E}[\mathbf{b}^\top \mathbf{b}].$$

By Eqn. (33),

$$(\boldsymbol{\theta}^*)^\top \mathbb{E}[\mathbf{A}^\top \mathbf{A}] \boldsymbol{\theta}^* \leq 2(1+\gamma) \|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi}^2.$$

And by Lemma G.4,

$$\mathbb{E}[\mathbf{b}^\top \mathbf{b}] \leq 9K(1-\gamma)^2.$$

To summarize,

$$\text{tr}(\boldsymbol{\Sigma}_e) \leq 4(1+\gamma) \|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi}^2 + 18K(1-\gamma)^2.$$

□

### E.3 Proof of Lemma 5.3

*Proof.* For simplicity, we use the same abbreviations as in Appendix E.2. As in the proof of [Lemma 2 54], we only need to show that, for any  $p \in \mathbb{N}$ ,  $\alpha \in (0, (1 - \sqrt{\gamma})/(38p))$ , it holds that

$$\mathbb{E} \left\{ \left[ (\mathbf{I}_{dK} - \alpha \mathbf{A})^\top (\mathbf{I}_{dK} - \alpha \mathbf{A}) \right]^p \right\} \leq \mathbf{I}_{dK} - \frac{1}{2} \alpha p (1 - \gamma) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \left( \leq \left( 1 - \frac{1}{2} \alpha p (1 - \gamma) \lambda_{\min} \right) \mathbf{I}_{dK} \right).$$

Let  $\mathbf{B} := \mathbf{A} + \mathbf{A}^\top - \alpha \mathbf{A}^\top \mathbf{A}$  which satisfies  $(\mathbf{I}_{dK} - \alpha \mathbf{A})^\top (\mathbf{I}_{dK} - \alpha \mathbf{A}) = \mathbf{I}_{dK} - \alpha \mathbf{B}$ . To give an upper bound of  $\mathbb{E}[(\mathbf{I}_{dK} - \alpha \mathbf{B})^p]$ , it suffices to show that

$$\mathbb{E}[\mathbf{B}] \geq (1 - \sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi, \quad \mathbb{E}[\mathbf{B}^p] \leq \frac{17}{16} 4^p \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi, \quad \forall p \in \{2, 3, \dots\},$$

if we take  $\alpha \in (0, (1 - \sqrt{\gamma})/(2(1 + \gamma)))$ .

Given these results, we have, when  $\alpha \in (0, (1 - \sqrt{\gamma})/(38p))$ , it holds that

$$\begin{aligned} \mathbb{E}[(\mathbf{I}_{dK} - \alpha \mathbf{B})^p] &\leq \mathbf{I} - \alpha p \mathbb{E}[\mathbf{B}] + \sum_{l=2}^p \alpha^l \binom{p}{l} \mathbb{E}[\mathbf{B}^l] \\ &\leq \mathbf{I} - \left( \alpha p (1 - \sqrt{\gamma}) - \frac{17}{16} \sum_{l=2}^{\infty} (4\alpha p)^l \right) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \\ &= \mathbf{I} - \left( \alpha p (1 - \sqrt{\gamma}) - \frac{17\alpha^2 p^2}{1 - 4\alpha p} \right) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \\ &\leq \mathbf{I} - \frac{1}{2} \alpha p (1 - \sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \end{aligned}$$

**Lower Bound of  $\mathbb{E}[\mathbf{B}]$ .** To show  $\mathbb{E}[\mathbf{B}] \geq (1 - \sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi$ , we first show that  $\mathbb{E}[\mathbf{A} + \mathbf{A}^\top] \geq 2(1 - \sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi$ , which is equivalent to  $(\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \mathbb{E}[\mathbf{A} + \mathbf{A}^\top] (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \geq 2(1 - \sqrt{\gamma})$ , where

$$\mathbb{E}[\mathbf{A} + \mathbf{A}^\top] = 2(\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi) - \mathbb{E} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right] - \mathbb{E} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right]^\top.$$

Then, for any  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$  with  $\|\boldsymbol{\theta}\| = 1$ ,

$$\begin{aligned}
& \boldsymbol{\theta}^\top (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \mathbb{E} [\mathbf{A} + \mathbf{A}^\top] (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \boldsymbol{\theta} \\
&= 2 - 2\boldsymbol{\theta}^\top (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \mathbb{E} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right] (\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi)^{-\frac{1}{2}} \boldsymbol{\theta} \\
&\geq 2 - 2 \left\| \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes \left( \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \phi(s) \phi(s')^\top \boldsymbol{\Sigma}_\phi^{-\frac{1}{2}} \right) \right] \right\| \\
&\geq 2(1 - \sqrt{\gamma}),
\end{aligned}$$

where we used the result Eqn. (29).

Next, we give an upper bound for  $\mathbb{E} [\mathbf{A}^\top \mathbf{A}]$ , we need to compute the following terms: by Lemma A.4,

$$\begin{aligned}
[\mathbf{I}_K \otimes (\phi(s) \phi(s)^\top)]^2 &= \mathbf{I}_K \otimes (\phi(s) \phi(s)^\top \phi(s) \phi(s)^\top) \\
&= \|\phi(s)\|^2 \mathbf{I}_K \otimes (\phi(s) \phi(s)^\top) \\
&\leq \mathbf{I}_K \otimes (\phi(s) \phi(s)^\top).
\end{aligned} \tag{31}$$

Hence

$$\mathbb{E} [\mathbf{I}_K \otimes (\phi(s) \phi(s)^\top)]^2 \leq \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi.$$

And by Lemma A.4 and Lemma G.3,

$$\begin{aligned}
& \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right]^\top \left[ \left( \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s) \phi(s')^\top) \right] \\
&= \left( \mathbf{C}^{-\top} \tilde{\mathbf{G}}^\top(r) \mathbf{C}^\top \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s') \phi(s)^\top \phi(s) \phi(s')^\top) \\
&= \|\phi(s)\|^2 \left( \mathbf{C}^{-\top} \tilde{\mathbf{G}}^\top(r) \mathbf{C}^\top \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right) \otimes (\phi(s') \phi(s')^\top) \\
&\leq \left\| \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right\|^2 \mathbf{I}_K \otimes (\phi(s') \phi(s')^\top) \\
&\leq \gamma \mathbf{I}_K \otimes (\phi(s') \phi(s')^\top).
\end{aligned}$$

To summarize, by the basic inequality  $(\mathbf{B}_1 - \mathbf{B}_2)^\top (\mathbf{B}_1 - \mathbf{B}_2) \leq 2(\mathbf{B}_1^\top \mathbf{B}_1 + \mathbf{B}_2^\top \mathbf{B}_2)$ , we have

$$\mathbf{A}^\top \mathbf{A} \leq 2\mathbf{I}_K \otimes (\phi(s) \phi(s)^\top + \gamma \phi(s') \phi(s')^\top), \tag{32}$$

and, after taking expectation,

$$\mathbb{E} [\mathbf{A}^\top \mathbf{A}] \leq 2(1 + \gamma) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi. \tag{33}$$

Combining these together, we obtain

$$\mathbb{E}[\mathbf{B}] \geq 2[(1 - \sqrt{\gamma}) - \alpha(1 + \gamma)] \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \geq (1 - \sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi,$$

if we take  $\alpha \in (0, (1 - \sqrt{\gamma})/(2(1 + \gamma)))$ .

**Upper Bound of  $\mathbb{E}[\mathbf{B}^p]$ .** Because  $\mathbf{B}^2$  is always PSD, we have the following upper bound

$$\mathbf{B}^p \leq \|\mathbf{B}\|^{p-2} \mathbf{B}^2.$$

We first give an almost-sure upper bound for  $\|\mathbf{B}\|$ . By Lemma 5.2,  $\|\mathbf{A}\| \leq 1 + \sqrt{\gamma}$ . And by Eqn. (32),

$$\begin{aligned} \|\mathbf{A}^\top \mathbf{A}\| &\leq 2 \|\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top + \gamma\phi(s')\phi(s')^\top)\| \\ &\leq 2 \|\phi(s)\phi(s)^\top + \gamma\phi(s')\phi(s')^\top\| \\ &\leq 2(1 + \gamma). \end{aligned} \tag{34}$$

Hence,

$$\begin{aligned} \|\mathbf{B}\| &= \|\mathbf{A} + \mathbf{A}^\top - \alpha\mathbf{A}^\top \mathbf{A}\| \\ &\leq 2\|\mathbf{A}\| + \alpha\|\mathbf{A}^\top \mathbf{A}\| \\ &\leq 2(1 + \sqrt{\gamma}) + 2\alpha(1 + \gamma) \\ &\leq 4, \end{aligned} \tag{35}$$

because  $\alpha \in (0, (1 - \sqrt{\gamma})/(2(1 + \gamma)))$ .

Now, we aim to give an upper bound for  $\mathbb{E}[\mathbf{B}^2]$ ,

$$\begin{aligned} \mathbf{B}^2 &= (\mathbf{A} + \mathbf{A}^\top - \alpha\mathbf{A}^\top \mathbf{A})^2 \\ &\leq (1 + \beta) (\mathbf{A} + \mathbf{A}^\top)^2 + (1 + \beta^{-1}) \alpha^2 (\mathbf{A}^\top \mathbf{A})^2 \\ &\leq 2(1 + \beta) (\mathbf{A}^\top \mathbf{A} + \mathbf{A} \mathbf{A}^\top) + (1 + \beta^{-1}) \alpha^2 (\mathbf{A}^\top \mathbf{A})^2, \end{aligned}$$

where we used the fact that  $(\mathbf{B}_1 + \mathbf{B}_2)^2 \leq (1 + \beta)\mathbf{B}_1^2 + (1 + \beta^{-1})\mathbf{B}_2^2$  for any symmetric matrices  $\mathbf{B}_1, \mathbf{B}_2$ , since  $\beta\mathbf{B}_1^2 + \beta^{-1}\mathbf{B}_2^2 - \mathbf{B}_1\mathbf{B}_2 - \mathbf{B}_2\mathbf{B}_1 = \left(\sqrt{\beta}\mathbf{B}_1 - \sqrt{\beta^{-1}}\mathbf{B}_2\right)^2 \geq \mathbf{0}$ ,  $\beta \in (0, 1)$  to be determined; and the fact that  $\mathbf{A}^2 + (\mathbf{A}^\top)^2 \leq \mathbf{A}^\top \mathbf{A} + \mathbf{A} \mathbf{A}^\top$  since the square of the skew-symmetric

matrix is negative semi-definite  $(\mathbf{A} - \mathbf{A}^\top)^2 \preceq \mathbf{0}$ . By Eqn. (34) and Eqn. (33), we have

$$\|\mathbf{A}^\top \mathbf{A}\| \leq 2(1 + \gamma),$$

$$\mathbb{E}[\mathbf{A}^\top \mathbf{A}] \preceq 2(1 + \gamma) \mathbf{I}_K \otimes \Sigma_\phi, \quad (36)$$

thus, by  $\alpha \in (0, (1 - \sqrt{\gamma})/(2(1 + \gamma)))$ , it holds that

$$\alpha^2 \mathbb{E}[(\mathbf{A}^\top \mathbf{A})^2] \leq 4\alpha^2(1 + \gamma)^2 \mathbf{I}_K \otimes \Sigma_\phi \preceq (1 - \sqrt{\gamma})^2 \mathbf{I}_K \otimes \Sigma_\phi. \quad (37)$$

As for  $\mathbb{E}[\mathbf{A}\mathbf{A}^\top]$ , by the basic inequality  $(\mathbf{B}_1 - \mathbf{B}_2)(\mathbf{B}_1 - \mathbf{B}_2)^\top \preceq 2(\mathbf{B}_1\mathbf{B}_1^\top + \mathbf{B}_2\mathbf{B}_2^\top)$ , we have

$$\begin{aligned} \mathbf{A}\mathbf{A}^\top &= \left\{ [\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)] - \left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right] \right\} \\ &\quad \cdot \left\{ [\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)] - \left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right] \right\}^\top \\ &\preceq 2[\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)]^2 + 2\left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right] \left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right]^\top. \end{aligned}$$

By Eqn. (31), we have

$$[\mathbf{I}_K \otimes (\phi(s)\phi(s)^\top)]^2 \preceq \mathbf{I}_K \otimes (\phi(s)\phi(s)^\top).$$

And by Lemma G.3,

$$\begin{aligned} &\left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right] \left[ (C\tilde{G}(r)C^{-1}) \otimes (\phi(s)\phi(s')^\top) \right]^\top \\ &= (C\tilde{G}(r)C^{-1}C^{-T}\tilde{G}^\top(r)C^\top) \otimes (\phi(s)\phi(s')^\top\phi(s')\phi(s)^\top) \\ &= \|\phi(s')\|^2 (C\tilde{G}(r)C^{-1}C^{-T}\tilde{G}^\top(r)C^\top) \otimes (\phi(s)\phi(s)^\top) \\ &\leq \|C\tilde{G}(r)C^{-1}\|^2 \mathbf{I}_K \otimes (\phi(s)\phi(s)^\top) \\ &\leq \gamma \mathbf{I}_K \otimes (\phi(s)\phi(s)^\top). \end{aligned}$$

To summarize, we have

$$\mathbf{A}\mathbf{A}^\top \preceq 2(1 + \gamma) \mathbf{I}_K \otimes (\phi(s)\phi(s)^\top),$$

and after taking expectation

$$\mathbb{E}[\mathbf{A}\mathbf{A}^\top] \preceq 2(1 + \gamma) \mathbf{I}_K \otimes \Sigma_\phi. \quad (38)$$



By putting everything together (Eqn. (36), Eqn. (37), Eqn. (38)), we have

$$\begin{aligned}
\mathbf{B}^2 &\leq 2(1 + \beta) (\mathbf{A}^\top \mathbf{A} + \mathbf{A} \mathbf{A}^\top) + (1 + \beta^{-1}) \alpha^2 (\mathbf{A}^\top \mathbf{A})^2 \\
&\leq [8(1 + \gamma)(1 + \beta) + (1 + \beta^{-1})(1 - \sqrt{\gamma})^2] \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \\
&\leq (17 + 9\gamma - 10\sqrt{\gamma}) \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \\
&\leq 17 \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi,
\end{aligned}$$

where we take  $\beta = \frac{1 - \sqrt{\gamma}}{\sqrt{8(1 + \gamma)}}$ . Therefore, by Eqn. (35),

$$\begin{aligned}
\mathbf{B}^p &\leq \|\mathbf{B}\|^{p-2} \mathbf{B}^2 \\
&\leq 4^{p-2} 17 \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi \\
&= \frac{17}{16} 4^p \mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi.
\end{aligned}$$

□

#### E.4 $L^p$ Convergence

**Theorem E.1** ( $L^p$  Convergence). *For any  $K \geq (1 - \gamma)^{-1}$ ,  $p > 2$  and  $\alpha \in (0, (1 - \sqrt{\gamma})/[38(p + \log T)])$ , it holds that*

$$\begin{aligned}
\mathbb{E}^{1/p} [(\mathcal{L}(\bar{\boldsymbol{\theta}}_T))^p] &\lesssim \frac{\sqrt{p}}{\sqrt{T}} \frac{\frac{1}{\sqrt{K(1-\gamma)}} \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + 1}{(1 - \gamma)\sqrt{\lambda_{\min}}} \left( 1 + \frac{\sqrt{\alpha p} + \alpha p}{\sqrt{(1 - \gamma)\lambda_{\min}}} \right) \\
&\quad + \frac{p}{T} \frac{\frac{1}{\sqrt{K(1-\gamma)}} \|\boldsymbol{\theta}^\star\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + 1}{(1 - \gamma)^{\frac{3}{2}} \lambda_{\min}} \left( 1 + \frac{1}{\sqrt{\alpha p}} \right) \\
&\quad + \frac{1}{T} \frac{(1 - \frac{1}{2}\alpha(1 - \sqrt{\gamma})\lambda_{\min})^{T/2}}{\sqrt{\alpha}(1 - \gamma)\sqrt{\lambda_{\min}}} \left( \frac{1}{\sqrt{\alpha}} + \frac{p}{\sqrt{(1 - \gamma)\lambda_{\min}}} \right) \frac{1}{\sqrt{K(1 - \gamma)}} \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^\star\|.
\end{aligned}$$

*Proof.* Combining Lemma 5.1, Lemma 5.2 and Lemma 5.3 with [Theorem 2 54], we have

$$\begin{aligned}
& \mathbb{E}^{1/p} [(\mathcal{L}(\bar{\boldsymbol{\theta}}_T))^p] \\
& \lesssim \frac{1}{\sqrt{K(1-\gamma)^4\lambda_{\min}}} \left[ \sqrt{\frac{p \operatorname{tr}(\boldsymbol{\Sigma}_e)}{T}} \left( 1 + \frac{C_A \sqrt{\alpha p}}{\sqrt{a}} + \frac{C_A C_e \alpha p}{\sqrt{\operatorname{tr}(\boldsymbol{\Sigma}_e)}} \right) + \frac{(1+C_A)C_e p}{T} \right. \\
& \quad \left. + \frac{p \sqrt{\operatorname{tr}(\boldsymbol{\Sigma}_e)}}{\sqrt{a}T} \left( 1 + \frac{1}{\sqrt{\alpha p}} \right) + (1-\alpha a)^{T/2} \left( \frac{1}{\alpha T} + \frac{C_{Ap}}{\sqrt{\alpha a}T} \right) \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\| \right] \\
& \lesssim \frac{1}{\sqrt{K(1-\gamma)^4\lambda_{\min}}} \left[ \sqrt{p} \frac{\|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + \sqrt{K}(1-\gamma)}{\sqrt{T}} \left( 1 + \frac{\sqrt{\alpha p} + \alpha p}{\sqrt{(1-\gamma)\lambda_{\min}}} \right) + \frac{p(\|\boldsymbol{\theta}^*\| + \sqrt{K}(1-\gamma))}{T} \right. \\
& \quad \left. + p \frac{\|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + \sqrt{K}(1-\gamma)}{\sqrt{(1-\gamma)\lambda_{\min}}T} \left( 1 + \frac{1}{\sqrt{\alpha p}} \right) \right. \\
& \quad \left. + (1 - \frac{1}{2}\alpha(1-\sqrt{\gamma})\lambda_{\min})^{T/2} \left( \frac{1}{\alpha T} + \frac{p}{\sqrt{\alpha(1-\gamma)\lambda_{\min}}T} \right) \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\| \right] \\
& \lesssim \frac{\sqrt{p}}{\sqrt{T}} \frac{1}{(1-\gamma)\sqrt{\lambda_{\min}}} \frac{\|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + 1}{\sqrt{K}(1-\gamma)} \left( 1 + \frac{\sqrt{\alpha p} + \alpha p}{\sqrt{(1-\gamma)\lambda_{\min}}} \right) \\
& \quad + \frac{p}{T} \frac{1}{(1-\gamma)^{\frac{3}{2}}\lambda_{\min}} \frac{\|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi} + 1}{\sqrt{K}(1-\gamma)} \left( 1 + \frac{1}{\sqrt{\alpha p}} \right) \\
& \quad + \frac{1}{T} \frac{(1 - \frac{1}{2}\alpha(1-\sqrt{\gamma})\lambda_{\min})^{T/2}}{\sqrt{\alpha(1-\gamma)}\sqrt{\lambda_{\min}}} \left( \frac{1}{\sqrt{\alpha}} + \frac{p}{\sqrt{(1-\gamma)\lambda_{\min}}} \right) \frac{1}{\sqrt{K}(1-\gamma)} \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|,
\end{aligned}$$

where we used the fact that  $\|\boldsymbol{\theta}^*\| \leq (\lambda_{\min})^{-1/2} \|\boldsymbol{\theta}^*\|_{\mathbf{I}_K \otimes \boldsymbol{\Sigma}_\phi}$ .  $\square$

## E.5 Convergence Results for SSGD with the PMF Representation

In this section, we present the counterparts of Lemma 5.1, Lemma 5.2, Lemma 5.3 and Theorem 4.1 for stochastic semi-gradient descent (SSGD) with the probability mass function (PMF) representation. These results will additionally depend on  $K$ . The additional  $K$ -dependent terms arise because the condition number of  $\mathbf{C}^\top \mathbf{C}$  scales with  $K^2$  (Lemma G.2). These terms are inevitable within our theoretical framework. The proofs of these results require only minor modifications to the original proofs, and we omit them for brevity.

In fact, in Appendix F, we validate some theoretical results through numerical experiments. To be concrete, we find that empirically, as  $K$  increases, to ensure convergence, the step size of the vanilla algorithm in [4, Section 9.6] indeed needs to decay at a rate of  $K^{-2}$ . In contrast, the step size of our **Linear-CTD** does not need to be adjusted when  $K$  increases. Moreover, we find that **Linear-CTD** empirically consistently outperforms the vanilla algorithm under different  $K$ .

Recall Eqn. (26), the updating scheme of the algorithm is

$$\begin{aligned}\mathbf{W}_t &\leftarrow \mathbf{W}_{t-1} - \alpha \phi(s_t) \left( \mathbf{p}_{\mathbf{w}_{t-1}}(s_t) - \mathbf{p}_{\mathcal{T}_t^\pi} \eta_{\mathbf{w}_{t-1}}(s_t) \right)^\top \mathbf{C}^\top \mathbf{C} \\ &= \mathbf{W}_{t-1} - \alpha \phi(s_t) \left[ \phi(s_t)^\top \mathbf{W}_{t-1} - \phi(s_{t+1})^\top \mathbf{W}_{t-1} \tilde{\mathbf{G}}^\top(r_t) - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^\top \right] \mathbf{C}^\top \mathbf{C},\end{aligned}$$

which is equivalent to

$$\mathbf{W}_t \mathbf{C}^\top \leftarrow \mathbf{W}_{t-1} \mathbf{C}^\top - \alpha \phi(s_t) \left[ \phi(s_t)^\top \mathbf{W}_{t-1} \mathbf{C}^\top - \phi(s_{t+1})^\top \mathbf{W}_{t-1} \mathbf{C}^\top (\mathbf{C} \tilde{\mathbf{G}}(r_t) \mathbf{C}^{-1})^\top - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^\top \mathbf{C}^\top \right] \mathbf{C} \mathbf{C}^\top,$$

here we drop the additional  $2\iota_K$  in the step size for brevity. Letting  $\boldsymbol{\Theta}_{\text{PMF},t} := \mathbf{W}_t \mathbf{C}^\top$  be the CDF parameter, the algorithm becomes

$$\boldsymbol{\Theta}_{\text{PMF},t} \leftarrow \boldsymbol{\Theta}_{\text{PMF},t-1} - \alpha \phi(s_t) \left[ \phi(s_t)^\top \boldsymbol{\Theta}_{\text{PMF},t-1} - \phi(s_{t+1})^\top \boldsymbol{\Theta}_{\text{PMF},t-1} (\mathbf{C} \tilde{\mathbf{G}}(r_t) \mathbf{C}^{-1})^\top - \frac{1}{K+1} \left( \sum_{j=0}^K \mathbf{g}_j(r_t) - \mathbf{1}_K \right)^\top \mathbf{C}^\top \right] \mathbf{C} \mathbf{C}^\top. \quad (39)$$

Here, we add the subscript PMF to the original notations to indicate the difference. Then, the algorithm corresponds to the following linear system for  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$

$$\bar{\mathbf{A}}_{\text{PMF}} \boldsymbol{\theta} = \bar{\mathbf{b}}_{\text{PMF}},$$

where

$$\begin{aligned}\bar{\mathbf{A}}_{\text{PMF}} &= [(\mathbf{C} \mathbf{C}^\top) \otimes \boldsymbol{\Sigma}_\phi] - \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C} \mathbf{C}^\top (\mathbf{C} \tilde{\mathbf{G}}(r_t) \mathbf{C}^{-1}) \right) \otimes (\phi(s) \phi(s')^\top) \right], \\ \bar{\mathbf{b}}_{\text{PMF}} &= \frac{1}{K+1} \mathbb{E}_{s,r} \left\{ \left[ \mathbf{C} \mathbf{C}^\top \mathbf{C} \left( \sum_{j=0}^K \mathbf{g}_j(r) - \mathbf{1}_K \right) \right] \otimes \phi(s) \right\}.\end{aligned}$$

Compared to our **Linear-CTD** (Eqn. (16)), this algorithm has an additional matrix  $\mathbf{C} \mathbf{C}^\top$  with the condition number of order  $K^2$  (see Lemma G.2).

Now, we are ready to state the theoretical results for the algorithm.

**Lemma E.1.** *For any  $\boldsymbol{\theta} \in \mathbb{R}^{dK}$ , it holds that*

$$\mathcal{L}(\boldsymbol{\theta}) \lesssim \frac{1}{\sqrt{K}(1-\gamma)^2 \sqrt{\lambda_{\min}}} \|\bar{\mathbf{A}}_{\text{PMF}} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|. \quad (40)$$

Lemma E.1 achieves the same order of bound as prior results for **Linear-CTD** (Lemma 5.1), as the minimum eigenvalue of  $\mathbf{C} \mathbf{C}^\top$  remains  $\Omega(1)$  (Lemma G.2). However, from the numerical experiments

(Figure 5) in Appendix F.2, we observe that after substituting  $\bar{\boldsymbol{\theta}}_t$  into  $\boldsymbol{\theta}$ , as  $K$  grows, the RHS grows with  $K$ , while the LHS remains almost unchanged. This might be because when the matrix  $\mathbf{C}\mathbf{C}^\top$  acts on the relevant random vectors, the stretching coefficient (*i.e.*,  $\|\mathbf{C}\mathbf{C}^\top \mathbf{x}\|/\|\mathbf{x}\|$  for some vector  $\mathbf{x}$ ) is usually of order  $K^2$  rather than a constant order. For example, consider the case where the matrix  $\mathbf{C}\mathbf{C}^\top$  acts on a random vector  $\mathbf{X}$  that follows a uniform distribution over the surface of unit sphere ( $\|\mathbf{X}\| = 1$ ). Since the  $k$ -th largest eigenvalue of the matrix  $\mathbf{C}\mathbf{C}^\top$  is of order  $k^2$ , by Hanson-Wright inequality [62, Theorem 6.2.1], we have  $\|\mathbf{C}\mathbf{C}^\top \mathbf{X}\|$  is of order  $K^2$  with high probability.

**Lemma E.2.** *It holds that*

$$C_{A,\text{PMF}} \lesssim K^2 C_A, \quad C_{e,\text{PMF}} \lesssim K^2 C_e, \quad \text{tr}(\boldsymbol{\Sigma}_{e,\text{PMF}}) \lesssim K^4 \text{tr}(\boldsymbol{\Sigma}_e).$$

Lemma E.2 introduces an additional factor of  $K^2$  (or  $K^4$ ) compared to previous results for Linear-CTD (Lemma 5.2), since the maximum eigenvalue of  $\mathbf{C}\mathbf{C}^\top$  is of order  $K^2$ .

**Lemma E.3.** *For any  $p \geq 2$ , let  $a_{\text{PMF}} \simeq (1 - \sqrt{\gamma})\lambda_{\min}$  and  $\alpha_{p,\infty}^{\text{PMF}} \simeq (1 - \sqrt{\gamma})/(pK^2)$  ( $\alpha_{p,\infty}^{\text{PMF}} p \leq 1/2$ ). Then for any  $\alpha \in (0, \alpha_{p,\infty}^{\text{PMF}})$ ,  $\mathbf{u} \in \mathbb{R}^{dK}$  and  $t \in \mathbb{N}$*

$$\mathbb{E}^{1/p} \left[ \left\| \boldsymbol{\Gamma}_{t,\text{PMF}}^{(\alpha)} \mathbf{u} \right\|^p \right] \leq (1 - \alpha a_{\text{PMF}})^t \|\mathbf{u}\|.$$

As before, in this lemma,  $a_{\text{PMF}}$  does not depend on  $K$  because the minimum eigenvalue of  $\mathbf{C}\mathbf{C}^\top$  is  $\Omega(1)$ , and  $\alpha_{p,\infty}^{\text{PMF}}$  scales with  $K^{-2}$  because the maximum eigenvalue of  $\mathbf{C}\mathbf{C}^\top$  is of order  $K^2$ .

**Theorem E.2.** *For any  $K \geq (1 - \gamma)^{-1}$  and  $\alpha \in (0, \alpha_{p,\infty}^{\text{PMF}})$ , it holds that*

$$\begin{aligned} \mathbb{E}^{1/2}[(\mathcal{L}(\bar{\boldsymbol{\theta}}_{\text{PMF},T}))^2] &\lesssim \frac{K^2 (\|\boldsymbol{\theta}^*\|_{V_1} + 1)}{\sqrt{T}(1 - \gamma)\sqrt{\lambda_{\min}}} \left( 1 + K \sqrt{\frac{\alpha K^2}{(1 - \gamma)\lambda_{\min}}} \right) + \frac{K^3 (\|\boldsymbol{\theta}^*\|_{V_1} + 1)}{T \sqrt{\alpha K^2} (1 - \gamma)^{\frac{3}{2}} \lambda_{\min}} \\ &\quad + K \frac{(1 - \frac{1}{2}(\alpha K^2) K^{-2} (1 - \sqrt{\gamma}) \lambda_{\min})^{T/2}}{T \sqrt{\alpha K^2} (1 - \gamma) \sqrt{\lambda_{\min}}} \left( \frac{K}{\sqrt{\alpha K^2}} + \frac{K^2}{\sqrt{(1 - \gamma)\lambda_{\min}}} \right) \|\boldsymbol{\theta}_{\text{PMF},0} - \boldsymbol{\theta}^*\|_{V_2}. \end{aligned}$$

This theorem for the PMF version algorithm yields an upper bound that is  $K^3$  times looser than Theorem 4.1 for our Linear-CTD. The appearance of the  $K^3$  factor is due to the fact that the condition number of the redundant matrix  $\mathbf{C}\mathbf{C}^\top$  is of order  $K^2$ . This factor is unavoidable within our theoretical analysis framework.

However, from the numerical experiments (Table 1 and Table 2) in Appendix F.2, we can only observe that our Linear-CTD consistently outperforms the PMF version algorithm under different

values of  $K$ , but the performance gap does not increase significantly when  $K$  becomes larger as predicted by Theorem 4.1 and Theorem E.2. The reason for this might be, as discussed after Lemma E.1: in the experimental environment we have set, when the matrix  $\mathbf{C}\mathbf{C}^\top$  acts on the vectors it encountered, the stretching coefficient is usually of order  $K^2$  rather than a constant order. See the numerical experiments (Figure 5) in Appendix F.2 for some evidence.

## F Numerical Experiment

In this appendix, we validate the proposed **Linear-CTD** algorithm (Eqn. (13)) with numerical experiments, and show its advantage over the baseline algorithm, stochastic semi-gradient descent (SSGD) with the probability mass function (PMF) representation (Eqn. (39)).

To empirically evaluate our **Linear-CTD** algorithm, we consider a 3-state MDP with  $\gamma = 0.75$ . When the number of states is finite, we denote by  $\Phi = (\phi(s))_{s \in \mathcal{S}} \in \mathbb{R}^{d \times \mathcal{S}}$  the feature matrix. Here, we set the feature matrix  $\Phi$  to be a full-rank matrix in  $\mathbb{R}^{3 \times 3}$ . The following experiments share zero initialization  $\theta_0 = \mathbf{0}$  with `max_iteration=500000` and `batch_size=25`.

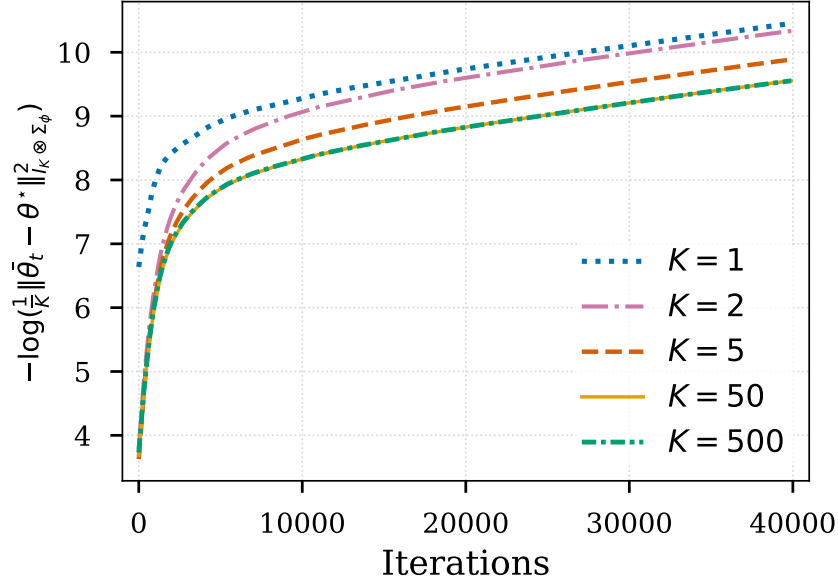
All of the experiments are conducted on a server with 4 NVIDIA RTX 4090 GPUs and Intel(R) Xeon(R) Gold 6132 CPU @ 2.60GHz.

### F.1 Empirical Convergence of Linear-CTD

We employ the **Linear-CTD** algorithm in the above environment and have the following convergence results in Figure 1. This figure shows the negative logarithm of  $\frac{1}{K} \|\bar{\theta}_t - \theta^\star\|_{I_K \otimes \Sigma_\phi}^2 = (1 - \gamma) \ell_{2, \mu_\pi}^2(\eta_\theta, \eta_{\theta^\star})$  along iterations. We observe that our **Linear-CTD** algorithm can converge for different values of  $K$  when we set the step size as  $\alpha = 0.01$ .

### F.2 Comparison with SSGD with the PMF Representation

First, we repeat the same experiment as in the previous section for the baseline algorithm, SSGD with the PMF representation. The experimental results in Figure 2 demonstrate that when the baseline algorithm uses a fixed step size  $\alpha = 0.01$ , it does not converge when  $K$  is large ( $K \geq 44$ ). The results in Figure 1 and Figure 2 verify the advantage of our **Linear-CTD** over the baseline algorithm as mentioned in Remark 5: when  $K$  increases, the step size of the baseline algorithm needs to decay (Lemma E.3). In contrast, the step size of our **Linear-CTD** does not need to be adjusted when  $K$  increases (Lemma 5.3).



**Figure 1.** Convergence results under varying  $K$  for our **Linear-CTD** algorithm with step size  $\alpha = 0.01$ . These curves exhibit similar trends, demonstrating our algorithm’s robustness across different  $K$  values.

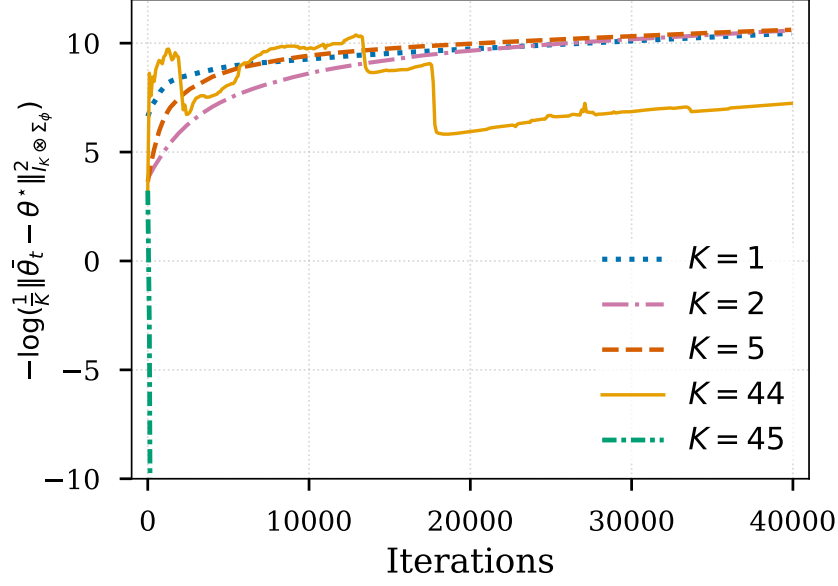
Next, we will verify that the maximum step size  $\alpha_\infty^{\text{PMF},(K)}$  that ensures the convergence of the baseline algorithm scales with  $K^{-2}$ , as predicted in Lemma E.3. Then we will compare the convergence rate of the baseline algorithm with that of our **Linear-CTD** algorithm.

In Figure 3, we employ the baseline algorithm with different step sizes under fixed  $K = 150$ , and we find that the baseline algorithm converges when the step size does not exceed  $8.6\text{e-}4$ , and it does not converge when the step size exceeds  $8.7\text{e-}4$ . This indicates that  $\alpha_\infty^{\text{PMF},(150)} \in [8.6\text{e-}4, 8.7\text{e-}4]$  in this environment, providing a good approximation of  $\alpha_\infty^{\text{PMF},(150)}$ .

We repeat the above experiments under varying  $K$ , searching for a step size that can ensure convergence (a lower bound of  $\alpha_\infty^{\text{PMF},(K)}$ ) and a step size that leads to divergence (an upper bound of  $\alpha_\infty^{\text{PMF},(K)}$ ) such that the two step sizes are as close as possible and thereby we can get a good approximation of  $\alpha_\infty^{\text{PMF},(K)}$ . The results are summarized in Table 1.

In Figure 4, we use the approximate values of  $\alpha_\infty^{\text{PMF},(K)}$  provided in Table 1 to perform a quadratic function fitting of  $1/\alpha_\infty^{\text{PMF},(K)}$  with respect to  $K$ . We find that  $\alpha_\infty^{\text{PMF},(K)}$  indeed approximately scales with  $K^{-2}$ , which verifies our theoretical result (Lemma E.3).

To compare the statistical efficiency of our **Linear-CTD** algorithm and the baseline algorithm, in Table 1, we also report the number of iterations required for the error to reach below  $2\text{e-}6$  when the step size satisfies  $\alpha \approx 0.2\alpha_\infty^{\text{PMF},(K)}$ . In addition, we present the parallel results of our **Linear-CTD** in Table 2. In Table 2, we find that the value of  $\alpha_\infty^{(K)}$  for our **Linear-CTD** algorithm is



**Figure 2.** Convergence results under varying  $K$  for the baseline algorithm, SSGD with the PMF representation with step size  $\alpha = 0.01$ . We remark that when  $K = 45$ , the program reports errors of inf and nan. In contrast to results of **Linear-CTD** in Figure 1, the baseline algorithm no longer converges when  $K$  is large ( $K \geq 44$ ).

much larger than  $\alpha_\infty^{\text{PMF},(K)}$ , and it does not decrease significantly with the growth of  $K$ . Moreover, by comparing the **Iterations** columns in Table 1 and Table 2, we find that the sample complexity of our **Linear-CTD** does not increase significantly with the growth of  $K$ , and **Linear-CTD** empirically consistently outperforms the baseline algorithm under different  $K$ .

However, the performance gap does not increase significantly as expected when  $K$  increases as predicted by Theorem 4.1 and Theorem E.2. The reason for this might be that, as discussed after Lemma E.1, in the experimental environment we have set, when the matrix  $\mathbf{C}\mathbf{C}^\top$  acts on the vectors it encountered, the stretching coefficient (*i.e.*,  $\|\mathbf{C}\mathbf{C}^\top \mathbf{x}\| / \|\mathbf{x}\|$  for some vector  $\mathbf{x}$ ) is usually of order  $K^2$  rather than a constant order.

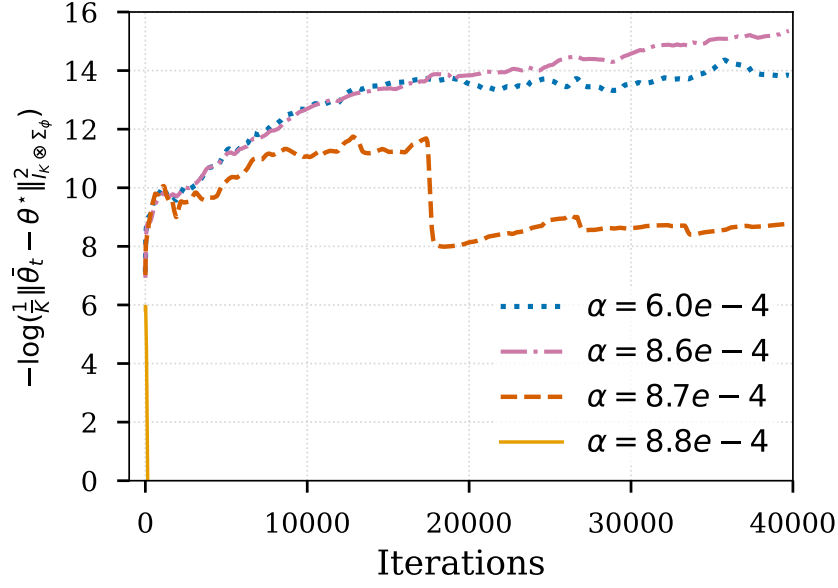
We verify this conjecture through the following experiment. We focus on the LHS and RHS of Eqn. (40) in Lemma E.1:

$$\mathcal{L}(\boldsymbol{\theta}) \lesssim \frac{1}{\sqrt{K}(1-\gamma)^2\sqrt{\lambda_{\min}}} \|\bar{\mathbf{A}}_{\text{PMF}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|,$$

where

$$\bar{\mathbf{A}}_{\text{PMF}} = [(\mathbf{C}\mathbf{C}^\top) \otimes \boldsymbol{\Sigma}_\phi] - \mathbb{E}_{s,r,s'} \left[ \left( \mathbf{C}\mathbf{C}^\top (\mathbf{C}\tilde{\mathbf{G}}(r_t)\mathbf{C}^{-1}) \right) \otimes (\phi(s)\phi(s')^\top) \right].$$

In our theoretical analysis, we first give an upper bound of the RHS, and then apply Lemma E.1



**Figure 3.** Convergence results with different step sizes for the baseline algorithm, SSGD with the PMF representation under fixed  $K = 150$ . We remark that when we take  $\alpha = 8.8e-4$ , the program reports errors of inf and nan. The baseline algorithm converges when the step size does not exceed  $8.6e-4$ , and it does not converge when the step size exceeds  $8.7e-4$ . Therefore,  $\alpha_{\infty}^{\text{PMF},(150)} \in [8.6e-4, 8.7e-4]$  in this environment.

bound the loss function  $\mathcal{L}(\theta)$  in the LHS with the RHS. However, since the minimum eigenvalue of the matrix  $\mathbf{C}\mathbf{C}^\top$  in the RHS is only of a constant order, we are unable to have a term of  $1/K^2$  in the RHS. Therefore, our conjecture can be verified by checking whether the bound provided in Lemma E.1 is tight in this environment, which is presented in Figure 5. The left sub-graph of Figure 5 corresponds to the LHS, and the right sub-graph corresponds to the RHS. We omit the constants that are independent of  $K$ . From Figure 5, we can find that the LHS remains almost unchanged under different  $K$ , but the RHS increases as  $K$  becomes larger. This indicates that the stretching coefficient of the matrix  $\mathbf{C}\mathbf{C}^\top$  that we frequently encounters during the iterative process grows with  $K$  rather than remaining a constant order. A similar analysis also holds for  $a_{\text{PMF}}$  in Lemma E.3, and we omit it for brevity. These factors result in the performance gap between our Linear-CTD algorithm and the baseline algorithm not increasing significantly when  $K$  becomes larger, as predicted by Theorem 4.1 and Theorem E.2.

## G Other Technical Lemmas

**Lemma G.1.** For any  $\nu_1, \nu_2 \in \mathcal{P}^{\text{sign}}$ , we have  $W_1(\nu_1, \nu_2) \leq \frac{1}{\sqrt{1-\gamma}} \ell_2(\nu_1, \nu_2)$ .



$K$	Lower Bound of $\alpha_{\infty}^{\text{PMF},(K)}$	Upper Bound of $\alpha_{\infty}^{\text{PMF},(K)}$	Iterations
30	2.1e-2	2.2e-2	37245
45	9e-3	9.5e-3	39262
75	3.4e-3	3.5e-3	38286
105	1.75e-3	1.8e-3	38123
150	8.6e-4	8.7e-4	38556
225	3.8e-4	3.9e-4	38317
300	2.1e-4	2.2e-4	38999
375	1.35e-4	1.4e-4	38674
450	9.5e-5	9.8e-5	38506

**Table 1.** Lower and upper bounds of the maximum step size  $\alpha_{\infty}^{\text{PMF},(K)}$  that ensures the convergence under varying  $K$  for the baseline algorithm, SSGD with the PMF representation. The bounds are determined using the same method as that in Figure 3. The **Iterations** column refers to the number of iterations required for the error to reach below  $2\text{e-}6$  when the step size satisfies  $\alpha \approx 0.2\alpha_{\infty}^{\text{PMF},(K)}$ .

$K$	Lower Bound of $\alpha_{\infty}^{(K)}$	Upper Bound of $\alpha_{\infty}^{(K)}$	Iterations
30	1.65	1.7	17908
45	1.65	1.7	17925
75	1.65	1.7	17942
105	1.6	1.65	18623
150	1.5	1.55	21947
225	1.55	1.6	20890
300	1.55	1.6	20890
375	1.55	1.65	20595
450	1.5	1.55	21947

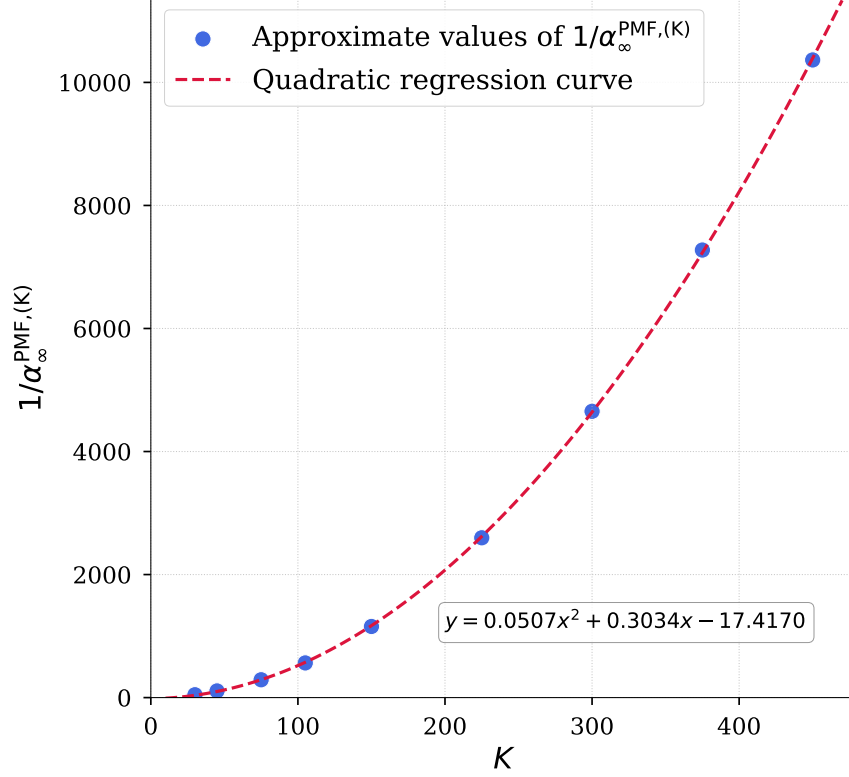
**Table 2.** Lower and upper bounds of the maximum step size  $\alpha_{\infty}^{(K)}$  that ensures the convergence under varying  $K$  for our **Linear-CTD**. The bounds are determined using the same method as that in Figure 3. The **Iterations** column refers to the number of iterations required for the error to reach below  $2\text{e-}6$  the step size satisfies  $\alpha \approx 0.2\alpha_{\infty}^{(K)}$ .

*Proof.* By Cauchy-Schwarz inequality,

$$\begin{aligned}
W_1(\nu_1, \nu_2) &= \int_0^{\frac{1}{1-\gamma}} |F_{\nu_1}(x) - F_{\nu_2}(x)| dx \\
&\leq \sqrt{\int_0^{\frac{1}{1-\gamma}} 1^2 dx} \sqrt{\int_0^{\frac{1}{1-\gamma}} |F_{\nu_1}(x) - F_{\nu_2}(x)|^2 dx} \\
&= \frac{1}{\sqrt{1-\gamma}} \ell_2(\nu_1, \nu_2).
\end{aligned}$$

□

**Lemma G.2.** Let  $C \in \mathbb{R}^{K \times K}$  be the matrix defined in Eqn. (11), it holds that the eigenvalues of



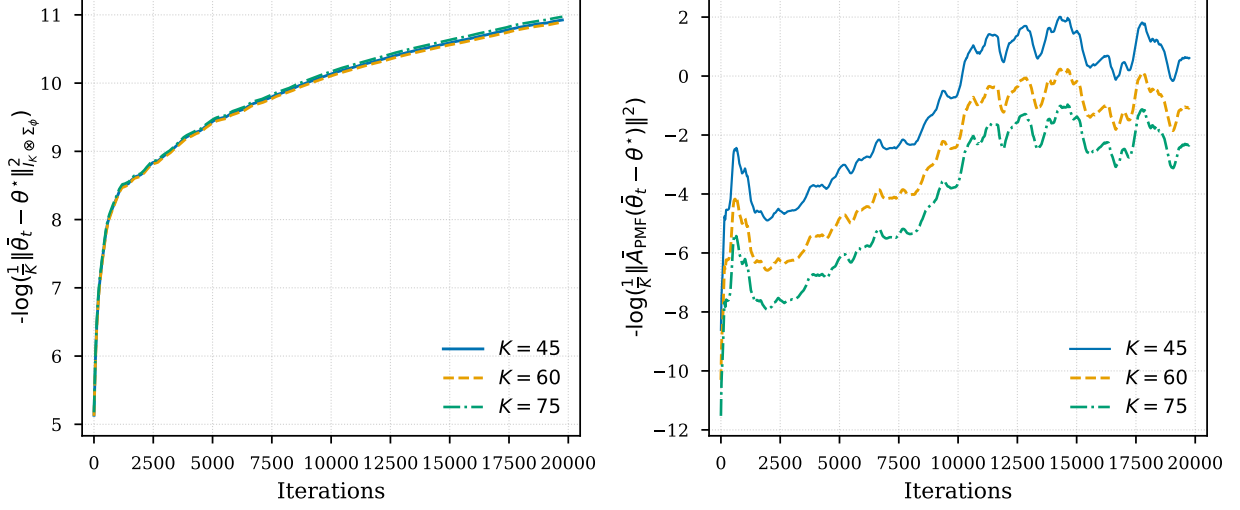
**Figure 4.** The approximate values of maximum step sizes  $1/\alpha_\infty^{\text{PMF},(K)}$  under varying  $K$ . Here we take the average of the upper and lower bounds of  $\alpha_\infty^{\text{PMF},(K)}$  provided in Table 1 as an approximation of  $\alpha_\infty^{\text{PMF},(K)}$  and perform quadratic regression of  $1/\alpha_\infty^{\text{PMF},(K)}$  on  $K$ . This fit achieves a mean squared error of 425.85 and  $R^2$  of 0.99996, which indicates that  $1/\alpha_\infty^{\text{PMF},(K)}$  indeed grows quadratically with respect to  $K$ , aligning with our theoretical results (Lemma E.3).

$C^T C$  are  $1/(4 \cos^2(k\pi/(2K+1)))$  for  $k \in [K]$ , and thus

$$\|C^T C\| = \frac{1}{4 \sin^2 \frac{\pi}{4K+2}} \leq K^2, \quad \|(C^T C)^{-1}\| = 4 \cos^2 \frac{\pi}{2K+1} \leq 4.$$

*Proof.* One can check that

$$C^T C = \begin{bmatrix} K & K-1 & \cdots & 2 & 1 \\ K-1 & K-1 & \cdots & 2 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 2 & 2 & \cdots & 2 & 1 \\ 1 & 1 & \cdots & 1 & 1 \end{bmatrix},$$



**Figure 5.** LHS and RHS of Eqn. (40) in Lemma E.1 under varying  $K$ . The left sub-graph corresponds to the LHS, and the right sub-graph corresponds to the RHS. We omit the constants that are independent of  $K$ . We can find that the LHS remains almost unchanged under different  $K$ , but the RHS increases as  $K$  becomes larger, indicating that the stretching coefficient of the matrix  $\mathbf{C}\mathbf{C}^\top$  that we frequently encounters during the iterative process grows with  $K$  rather than remaining a constant order.

$$(\mathbf{C}^\top \mathbf{C})^{-1} = \begin{bmatrix} 1 & -1 & 0 & \cdots & 0 & 0 \\ -1 & 2 & -1 & \cdots & 0 & 0 \\ 0 & -1 & 2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 2 & -1 \\ 0 & 0 & 0 & \cdots & -1 & 1 \end{bmatrix}.$$

Then, one can work with the the inverse of  $\mathbf{C}^\top \mathbf{C}$  and calculate its singular values by induction, which has similar forms to the analysis of Toeplitz's matrix. See Godsil [18] for more details.  $\square$

**Lemma G.3.** For any  $r \in [0, 1]$ , it holds that  $\|\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}\| \leq \sqrt{\gamma}$  and  $\|(\mathbf{C}^\top \mathbf{C})^{1/2} \tilde{\mathbf{G}}(r) (\mathbf{C}^\top \mathbf{C})^{-1/2}\| \leq \sqrt{\gamma}$ .

*Proof.* One can check that

$$\mathbf{C}^{-1} = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ -1 & 1 & \cdots & 0 & 0 \\ 0 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & -1 & 1 \end{bmatrix}.$$

It is clear that

$$\begin{aligned}
\left\| (\mathbf{C}^\top \mathbf{C})^{1/2} \tilde{\mathbf{G}}(r) (\mathbf{C}^\top \mathbf{C})^{-1/2} \right\| &\leq \sqrt{\gamma} \iff (\mathbf{C}^\top \mathbf{C})^{1/2} \tilde{\mathbf{G}}(r) (\mathbf{C}^\top \mathbf{C})^{-1} \tilde{\mathbf{G}}^\top(r) (\mathbf{C}^\top \mathbf{C})^{1/2} \preceq \gamma \mathbf{I}_K \\
&\iff \tilde{\mathbf{G}}(r) (\mathbf{C}^\top \mathbf{C})^{-1} \tilde{\mathbf{G}}^\top(r) \preceq \gamma (\mathbf{C}^\top \mathbf{C})^{-1} \\
&\iff \mathbf{C} \tilde{\mathbf{G}}(r) (\mathbf{C}^\top \mathbf{C})^{-1} \tilde{\mathbf{G}}^\top(r) \mathbf{C}^\top \preceq \gamma \mathbf{I}_K \\
&\iff \left\| \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right\| \leq \sqrt{\gamma}.
\end{aligned}$$

By Lemma H.2 and an upper bound on the spectral norm (Riesz–Thorin interpolation theorem) [55, Theorem 7.3], we obtain that

$$\left\| \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right\| \leq \sqrt{\left\| \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right\|_1 \left\| \mathbf{C} \tilde{\mathbf{G}}(r) \mathbf{C}^{-1} \right\|_\infty} \leq \sqrt{1 \cdot \gamma} = \sqrt{\gamma}.$$

□

**Lemma G.4.** Suppose  $K \geq (1 - \gamma)^{-1}$ ,  $\nu = (K + 1)^{-1} \sum_{k=0}^K \delta_{x_k}$  is the discrete uniform distribution, then for any  $r \in [0, 1]$ , it holds that

$$\ell_2((b_{r,\gamma})_\#(\nu), \nu) \leq 3\sqrt{1 - \gamma}.$$

*Proof.* Let  $\tilde{\nu}$  be the continuous uniform distribution on  $[0, (1 - \gamma)^{-1} + \iota_K]$ , we consider the following decomposition

$$\ell_2(\nu, (b_{r,\gamma})_\#(\nu)) \leq \ell_2(\nu, \tilde{\nu}) + \ell_2(\tilde{\nu}, (b_{r,\gamma})_\#(\tilde{\nu})) + \ell_2((b_{r,\gamma})_\#(\tilde{\nu}), (b_{r,\gamma})_\#(\nu)).$$

By definition, we have

$$\begin{aligned}
\ell_2(\nu, \tilde{\nu}) &= \sqrt{(K + 1) \int_0^{\iota_K} \left( (1 - \gamma) \frac{K}{K + 1} x \right)^2 dx} \\
&= \sqrt{\frac{1}{3K(K + 1)(1 - \gamma)}} \\
&\leq \frac{1}{K\sqrt{1 - \gamma}}.
\end{aligned}$$

By the contraction property, we have

$$\ell_2((b_{r,\gamma})_\#(\nu), (b_{r,\gamma})_\#(\tilde{\nu})) \leq \sqrt{\gamma} \ell_2(\nu, \tilde{\nu}) \leq \frac{\sqrt{\gamma}}{K\sqrt{1 - \gamma}}.$$

We only need to bound  $\ell_2(\tilde{\nu}, (b_{r,\gamma})_{\#}(\tilde{\nu}))$ . We can find that  $(b_{r,\gamma})_{\#}(\tilde{\nu})$  is the continuous uniform distribution on  $[r, r + \gamma\iota_K + \gamma(1 - \gamma)^{-1}]$ , and the upper bound is less than the upper bound of  $\nu$ , namely,  $r + \gamma\iota_K + \gamma(1 - \gamma)^{-1} \leq (1 - \gamma)^{-1} + \gamma\iota_K < (1 - \gamma)^{-1} + \iota_K$ . Hence

$$\begin{aligned} \ell_2^2(\tilde{\nu}, (b_{r,\gamma})_{\#}(\tilde{\nu})) &= \int_0^r \left( (1 - \gamma) \frac{K}{K + 1} x \right)^2 dx + \int_r^{r + \gamma\iota_K + \gamma(1 - \gamma)^{-1}} \left[ (1 - \gamma) \frac{K}{K + 1} \left( x - \frac{x - r}{\gamma} \right) \right]^2 dx \\ &\quad + \int_{r + \gamma\iota_K + \gamma(1 - \gamma)^{-1}}^{(1 - \gamma)^{-1} + \iota_K} \left( 1 - (1 - \gamma) \frac{K}{K + 1} x \right)^2 dx \\ &= \frac{(1 - \gamma)^2 K^2 r^3}{3(K + 1)^2} + \left( \frac{(1 - \gamma)\gamma K^2 r^3}{3(K + 1)^2} + \frac{(1 - \gamma)\gamma K^2 \left( \frac{K + 1}{K} - r \right)^3}{3(K + 1)^2} \right) + \frac{(1 - \gamma)^2 K^2 \left( \frac{K + 1}{K} - r \right)^3}{3(K + 1)^2} \\ &\leq (1 - \gamma)^2 + (1 - \gamma)\gamma \\ &= 1 - \gamma. \end{aligned}$$

To summarize, we have

$$\ell_2(\nu, (b_{r,\gamma})_{\#}(\nu)) \leq \frac{1}{K\sqrt{1 - \gamma}} + \sqrt{1 - \gamma} + \frac{\sqrt{\gamma}}{K\sqrt{1 - \gamma}} \leq 3\sqrt{1 - \gamma},$$

where we used the assumption  $K \geq (1 - \gamma)^{-1}$ . □

## H Analysis of the Categorical Projected Bellman Matrix

Recall that  $\tilde{\mathbf{G}}(r) = \mathbf{G}(r) - \mathbf{1}_K^\top \otimes \mathbf{g}_K(r)$ . We extend the definition in Theorem 3.1 and let  $g_{j,k}(r) = h((r + \gamma x_j - x_k)/\iota_K)_+ = h(r/\iota_K + \gamma j - k)$  for  $j, k \in \{0, 1, \dots, K\}$  where  $h(x) = (1 - |x|)_+$ .

**Lemma H.1.** *For any  $r \in [0, 1]$  and any  $k \in \{0, 1, \dots, K\}$ , in  $\mathbf{g}_k(r)$  there is either only one nonzero entry or two adjacent nonzero entries.*

*Proof.* It is clear that  $h(x) > 0 \iff -1 < x < 1$ . Let  $k_j(r) = \min \{k : g_{j,k}(r) > 0\}$ , then  $k_j(r) = \min \{k : r/\iota_K + \gamma j - k < 1\} = \min \{k : 0 \leq r/\iota_K + \gamma j - k < 1\}$ . The existence of  $k_j(r)$  is due to

$$r/\iota_K + \gamma j - K \leq 1/\iota_K + \gamma j - K \leq (1 - \gamma)K + \gamma K - K = 0 < 1.$$

Let  $a_j(r) := r/\iota_K + \gamma j - k_j(r) \in [0, 1]$ . Then  $g_{j,k_j(r)}(r) = h(a_j(r)) = 1 - a_j(r)$  and  $g_{j,k_j(r)+1}(r) = h(a_j(r) - 1) = a_j(r)$  are the only entries that can be nonzeros. □

The following results are immediate corollaries.

**Corollary H.1.**

$$\sum_{k=0}^i g_{j,k}(r) = \begin{cases} 0, & \text{for } 0 \leq i < k_j(r), \\ 1 - a_j(r), & \text{for } i = k_j(r), \\ 1, & \text{for } k_j(r) < i \leq K. \end{cases}$$

**Corollary H.2.**

$$k_{j+1}(r) = \begin{cases} k_j(r), & \text{if } a_j(r) \leq 1 - \gamma, \\ k_j(r) + 1, & \text{if } a_j(r) > 1 - \gamma. \end{cases}$$

As a result,

$$a_{j+1}(r) = \begin{cases} a_j(r) + \gamma, & \text{if } a_j(r) \leq 1 - \gamma, \\ a_j(r) + \gamma - 1, & \text{if } a_j(r) > 1 - \gamma. \end{cases}$$

**Lemma H.2.** All entries in  $\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}$  are non-negative.  $\|\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}\|_{\infty} = \gamma$  and  $\|\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}\|_1 \leq 1$ .

*Proof.* By definition the entries of  $\tilde{\mathbf{G}}(r)$  are

$$(\tilde{\mathbf{G}}(r))_{j,i} = g_{j,i}(r) - g_{K,i}(r) \quad \text{for } j, i \in \{0, 1, \dots, K-1\}.$$

Using the previous corollaries, through direct calculation we have that if  $k_{j+1}(r) = k_j(r)$ ,

$$(\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1})_{j,i} = \sum_{k=0}^i g_{j,k}(r) - \sum_{k=0}^i g_{j+1,k}(r) = \begin{cases} 0, & \text{for } 0 \leq i < k_j(r), \\ a_{j+1}(r) - a_j(r), & \text{for } i = k_j(r), \\ 0, & \text{for } k_j(r) < i < K. \end{cases}$$

And if  $k_{j+1}(r) = k_j(r) + 1$ ,

$$(\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1})_{j,i} = \sum_{k=0}^i g_{j,k}(r) - \sum_{k=0}^i g_{j+1,k}(r) = \begin{cases} 0, & \text{for } 0 < i < k_j(r), \\ 1 - a_j(r), & \text{for } i = k_j(r), \\ a_{j+1}(r), & \text{for } i = k_{j+1}(r), \\ 0, & \text{for } k_j(r) < i < K. \end{cases}$$

As a result, all entries in  $\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}$  is non-negative. Moreover, the sum of each column and

$\left\| \mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1} \right\|_{\infty}$  is  $\gamma$  since

$$\sum_{i=0}^{K-1} (\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1})_{j,i} = \begin{cases} a_{j+1}(r) - a_j(r) = \gamma, & \text{if } k_{j+1}(r) = k_j(r), \\ 1 - a_j(r) + a_{j+1}(r) = \gamma, & \text{if } k_{j+1}(r) = k_j(r) + 1. \end{cases}$$

Moreover, the row sum of  $\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1}$  is

$$\sum_{j=0}^{K-1} (\mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1})_{j,i} = \sum_{m=0}^i g_{0,m}(r) - \sum_{m=0}^i g_{K,m}(r) \leq 1 - 0 = 1.$$

Thus, it holds that  $\left\| \mathbf{C}\tilde{\mathbf{G}}(r)\mathbf{C}^{-1} \right\|_1 \leq 1$ . □